

Introduction to

Cosmology

Fourth Edition



Matts Roos

WILEY

*Introduction
to Cosmology*

Fourth Edition

Introduction to Cosmology

Fourth Edition

Matts Roos

WILEY

This edition first published 2015
© 2015 John Wiley & Sons, Ltd

Registered office

John Wiley & Sons Ltd, The Atrium, Southern Gate, Chichester, West Sussex, PO19 8SQ, United Kingdom

For details of our global editorial offices, for customer services and for information about how to apply for permission to reuse the copyright material in this book please see our website at www.wiley.com.

The right of the author to be identified as the author of this work has been asserted in accordance with the Copyright, Designs and Patents Act 1988.

All rights reserved. No part of this publication may be reproduced, stored in a retrieval system, or transmitted, in any form or by any means, electronic, mechanical, photocopying, recording or otherwise, except as permitted by the UK Copyright, Designs and Patents Act 1988, without the prior permission of the publisher.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic books.

Designations used by companies to distinguish their products are often claimed as trademarks. All brand names and product names used in this book are trade names, service marks, trademarks or registered trademarks of their respective owners. The publisher is not associated with any product or vendor mentioned in this book.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. It is sold on the understanding that the publisher is not engaged in rendering professional services and neither the publisher nor the author shall be liable for damages arising herefrom. If professional advice or other expert assistance is required, the services of a competent professional should be sought

Library of Congress Cataloging-in-Publication Data applied for.

A catalogue record for this book is available from the British Library.

ISBN 978-1-118-92332-0 (paperback)

Set in 9.5/12.5pt, NewAsterLTStd by Laserwords Private Limited, Chennai, India

To my family

Contents

Preface to First Edition	xi
Preface to Second Edition	xiii
Preface to Third Edition	xv
Preface to Fourth Edition	xvii
1 From Newton to Hubble	1
1.1 Historical Cosmology	2
1.2 Inertial Frames and the Cosmological Principle	6
1.3 Olbers' Paradox	8
1.4 Hubble's Law	11
1.5 The Age of the Universe	14
1.6 Matter in the Universe	16
1.7 Expansion in a Newtonian World	19
2 Special Relativity	25
2.1 Lorentz Transformations	25
2.2 Metrics of Curved Space-time	30
2.3 Relativistic Distance Measures	37
2.4 Tests of Special Relativity	45
3 General Relativity	49
3.1 The Principle of Equivalence	50
3.2 The Principle of Covariance	54
3.3 The Einstein Equation	58
3.4 Weak Field Limit	61

4	Tests of General Relativity	65
4.1	The Classical Tests	65
4.2	Binary Pulsars	67
4.3	Gravitational Lensing	69
4.4	Gravitational Waves	74
5	Cosmological Models	81
5.1	Friedmann–Lemaître Cosmologies	81
5.2	de Sitter Cosmology	93
5.3	The Schwarzschild Model	95
5.4	Black Holes	96
5.5	Extended Gravity Models	106
6	Thermal History of the Universe	111
6.1	Planck Time	112
6.2	The Primordial Hot Plasma	112
6.3	Electroweak Interactions	121
6.4	Photon and Lepton Decoupling	128
6.5	Big Bang Nucleosynthesis	134
6.6	Baryosynthesis and Antimatter Generation	142
7	Cosmic Inflation	151
7.1	Paradoxes of the Expansion	152
7.2	Consensus Inflation	158
7.3	The Chaotic Model	165
7.4	Predictions	168
7.5	A Cyclic Universe	169
8	Cosmic Microwave Background	175
8.1	The CMB Temperature	176
8.2	Temperature Anisotropies	180
8.3	Polarization	185
8.4	Model Testing and Parameter Estimation	189
9	Dark Matter	199
9.1	Virially Bound Systems	200
9.2	Galaxies	203
9.3	Clusters	208
9.4	Merging Galaxy Clusters	211
9.5	Dark Matter Candidates	213
9.6	The Cold Dark Matter Paradigm	218
10	Cosmic Structures	223
10.1	Density Fluctuations	223
10.2	Structure Formation	228

11 Dark Energy	235
11.1 The Cosmological Constant	235
11.2 Single Field Models	238
11.3 $f(R)$ Models	246
11.4 Extra Dimensions	248
12 Epilogue	255
Tables	257
Index	261

Preface to First Edition

A few decades ago, astronomy and particle physics started to merge in the common field of cosmology. The general public had always been more interested in the visible objects of astronomy than in invisible atoms, and probably met cosmology first in Steven Weinberg's famous book *The First Three Minutes*. More recently Stephen Hawking's *A Brief History of Time* has caused an avalanche of interest in this subject.

Although there are now many popular monographs on cosmology, there are so far no introductory textbooks at university undergraduate level. Chapters on cosmology can be found in introductory books on relativity or astronomy, but they cover only part of the subject. One reason may be that cosmology is explicitly cross-disciplinary, and therefore it does not occupy a prominent position in either physics or astronomy curricula.

At the University of Helsinki I decided to try to take advantage of the great interest in cosmology among the younger students, offering them a one-semester course about one year before their specialization started. Hence I could not count on much familiarity with quantum mechanics, general relativity, particle physics, astrophysics or statistical mechanics. At this level, there are courses with the generic name of Structure of Matter dealing with Lorentz transformations and the basic concepts of quantum mechanics. My course aimed at the same level. Its main constraint was that it had to be taught as a one-semester course, so that it would be accepted in physics and astronomy curricula. The present book is based on that course, given three times to physics and astronomy students in Helsinki.

Of course there already exist good books on cosmology. The reader will in fact find many references to such books, which have been an invaluable source of information to me. The problem is only that they address a postgraduate audience that intends to specialize in cosmology research. My readers will have to turn to these books later when they have mastered all the professional skills of physics and mathematics.

In this book I am not attempting to teach basic physics to astronomers. They will need much more. I am trying to teach just enough physics to be able to explain the

main ideas in cosmology without too much hand-waving. I have tried to avoid the other extreme, practised by some of my particle physics colleagues, of writing books on cosmology with the obvious intent of making particle physicists out of every theoretical astronomer.

I also do not attempt to teach basic astronomy to physicists. In contrast to astronomy scholars, I think the main ideas in cosmology do not require very detailed knowledge of astrophysics or observational techniques. Whole books have been written on distance measurements and the value of the Hubble parameter, which still remains imprecise to a factor of two. Physicists only need to know that quantities entering formulae are measurable—albeit incorporating factors h to some power—so that the laws can be discussed meaningfully. At undergraduate level, it is not even usual to give the errors on measured values.

In most chapters there are subjects demanding such a mastery of theoretical physics or astrophysics that the explanations have to be qualitative and the derivations meagre, for instance in general relativity, spontaneous symmetry breaking, inflation and galaxy formation. This is unavoidable because it just reflects the level of undergraduates. My intention is to go just a few steps further in these matters than do the popular monographs.

I am indebted in particular to two colleagues and friends who offered constructive criticism and made useful suggestions. The particle physicist Professor Kari Enqvist of NORDITA, Copenhagen, my former student, has gone to the trouble of reading the whole manuscript. The space astronomer Professor Stuart Bowyer of the University of California, Berkeley, has passed several early mornings of jet lag in Lapland going through the astronomy-related sections. Anyway, he could not go out skiing then because it was either a snow storm or -30°C ! Finally, the publisher provided me with a very knowledgeable and thorough referee, an astrophysicist no doubt, whose criticism of the chapter on galaxy formation was very valuable to me. For all remaining mistakes I take full responsibility. They may well have been introduced by me afterwards.

Thanks are also due to friends among the local experts: particle physicist Professor Masud Chaichian and astronomer Professor Kalevi Mattila have helped me with details and have answered my questions on several occasions. I am also indebted to several people who helped me to assemble the pictorial material: Drs Subir Sarkar in Oxford, Rocky Kolb in the Fermilab, Carlos Frenk in Durham, Werner Kienzle at CERN and members of the COBE team.

Finally, I must thank my wife Jacqueline for putting up with almost two years of near absence and full absent-mindedness while writing this book.

Matts Roos

Preface to Second Edition

In the three years since the first edition of this book was finalized, the field of cosmology has seen many important developments, mainly due to new observations with superior instruments such as the Hubble Space Telescope and the ground-based Keck telescope and many others. Thus a second edition has become necessary in order to provide students and other readers with a useful and up to date textbook and reference book.

At the same time I could balance the presentation with material which was not adequately covered before—there I am in debt to many readers. Also, the inevitable number of misprints, errors and unclear formulations, typical of a first edition, could be corrected. I am especially indebted to Kimmo Kainulainen who served as my course assistant one semester, and who worked through the book and the problems thoroughly, resulting in a very long list of corrigenda. A similar shorter list was also dressed by George Smoot and a student of his. It still worries me that the errors found by George had been found neither by Kimmo nor by myself, thus statistics tells me that some errors still will remain undetected.

For new pictorial material I am indebted to Wes Colley at Princeton, Carlos Frenk in Durham, Charles Lineweaver in Strasbourg, Jukka Nevalainen in Helsinki, Subir Sarkar in Oxford, and George Smoot in Berkeley. I am thankful to the Academie des Sciences for an invitation to Paris where I could visit the Observatory of Paris-Meudon and profit from discussions with S. Bonazzola and Brandon Carter.

Several of my students have contributed in various ways: by misunderstandings, indicating the need for better explanations, by their enthusiasm for the subject, and by technical help, in particular S. M. Harun-or-Rashid. My youngest grandchild Adrian (not yet 3) has showed a vivid interest for supernova bangs, as demonstrated by an X-ray image of the Cassiopeia A remnant. Thus the future of the subject is bright.

Matts Roos

Preface to Third Edition

This preface can start just like the previous one: in the seven years since the second edition was finalized, the field of cosmology has seen many important developments, mainly due to new observations with superior instruments. In the past, cosmology often relied on philosophical or aesthetic arguments; now it is maturing to become an exact science. For example, the Einstein–de Sitter universe, which has zero cosmological constant ($\Omega_\lambda = 0$), used to be favored for esthetical reasons, but today it is known to be very different from zero ($\Omega_\lambda = 0.73 \pm 0.04$).

In the first edition I quoted $\Omega_0 = 0.8 \pm 0.3$ (daring to believe in errors that many others did not), which gave room for all possible spatial geometries: spherical, flat and hyperbolic. Since then the value has converged to $\Omega_0 = 1.02 \pm 0.02$, and everybody is now willing to concede that the geometry of the Universe is flat, $\Omega_0 = 1$. This result is one of the cornerstones of what we now can call the ‘Concordance Model of Cosmology’. Still, deep problems remain, so deep that even Einstein’s general relativity is occasionally put in doubt.

A consequence of the successful march towards a ‘concordance model’ is that many alternative models can be discarded. An introductory text of limited length like the current one cannot be a historical record of failed models. Thus I no longer discuss, or discuss only briefly, $k \neq 0$ geometries, the Einstein–de Sitter universe, hot and warm dark matter, cold dark matter models with $\lambda = 0$, isocurvature fluctuations, topological defects (except monopoles), Bianchi universes, and formulae which only work in discarded or idealized models, like Mattig’s relation and the Saha equation.

Instead, this edition contains many new or considerably expanded subjects: Section 2.3 on Relativistic Distance Measures, Section 3.3 on Gravitational Lensing, Section 3.5 on Gravitational Waves, Section 4.3 on Dark Energy and Quintessence, Section 5.1 on Photon Polarization, Section 7.4 on The Inflaton as Quintessence, Section 7.5 on Cyclic Models, Section 8.3 on CMB Polarization Anisotropies, Section 8.4 on model testing and parameter estimation using mainly the first-year CMB results of the Wilkinson Microwave Anisotropy Probe, and Section 9.5

on large-scale structure results from the 2 degree Field (2dF) Galaxy Redshift Survey. The synopsis in this edition is also different and hopefully more logical, much has been entirely rewritten, and all parameter values have been updated.

I have not wanted to go into pure astrophysics, but the line between cosmology and cosmologically important astrophysics is not easy to draw. Supernova explosion mechanisms and black holes are included as in the earlier editions, but not for instance active galactic nuclei (AGNs) or jets or ultra-high-energy cosmic rays. Observational techniques are mentioned only briefly—they are beyond the scope of this book.

There are many new figures for which I am in debt to colleagues and friends, all acknowledged in the figure legends. I have profited from discussions with Professor Carlos Frenk at the University of Durham and Professor Kari Enqvist at the University of Helsinki. I am also indebted to Professor Juhani Keinonen at the University of Helsinki for having generously provided me with working space and access to all the facilities at the Department of Physical Sciences, despite the fact that I am retired.

Many critics, referees and other readers have made useful comments that I have tried to take into account. One careful reader, Urbana Lopes França Jr, sent me a long list of misprints and errors. A critic of the second edition stated that the errors in the first edition had been corrected, but that new errors had emerged in the new text. This will unfortunately always be true in any comparison of edition $n + 1$ with edition n . In an attempt to make continuous corrections I have assigned a web site for a list of errors and misprints.

My most valuable collaborator has been Thomas S. Coleman, a nonphysicist who contacted me after having spotted some errors in the second edition, and who proposed some improvements in case I were writing a third edition. This came at the appropriate time and led to a collaboration in which Thomas S. Coleman read the whole manuscript, corrected misprints, improved my English, checked my calculations, designed new figures and proposed clarifications where he found the text difficult.

My wife Jacqueline has many interesting subjects of conversation at the breakfast table. Regretfully, her breakfast companion is absent-minded, thinking only of cosmology. I thank her heartily for her kind patience, promising improvement.

Matts Roos

Preface to Fourth Edition

Just like the previous times I can state that the field of cosmology has seen so many important developments in the 11 years since the third edition that a fourth edition has become necessary.

In the first Chapter there is a new section presenting briefly various forms of baryonic matter: in supernovae and neutron stars and in noncollapsed objects such as interstellar dust, hot gas in the intergalactic medium, the cosmic rays, neutrinos and antiparticles. Active galactic nuclei, gamma ray bursts and quasars are also mentioned. The Hubble parameter and the age of the Universe are updated.

Chapter 2 is reorganized to contain only special relativity whereas all of general relativity forms Chapter 3. Chapter 2 ends with a brief new section on tests of special relativity and variable speed of light.

In the previous editions the Einstein equation was “derived” in the weak field limit with some hand-waving arguments. That derivation is still in Chapter 3, but the Einstein equation is now properly derived from the Hilbert–Einstein action. To some readers the derivation will then be more difficult, but one can of course skip it and be satisfied with the weak field limit. How to incorporate the energy of a gravitational field as a pseudotensor is briefly mentioned.

Chapter 4 addresses tests of general relativity with the exception for black holes which are now in Chapter 5. A new binary pulsar has been added and the section on detectors of gravitational radiation has been updated.

Chapter 5 on cosmological models contains the Schwarzschild model (previously in Chapter 2), a considerably expanded and modernized section on black holes, somewhat speculative perhaps, because it is based on Hawking’s recent ideas. The last section is entirely new, it discusses extended gravity models starting from a generalization of the Einstein–Hilbert action, and it lays the basis for dark energy models in Chapter 11.

Chapter 6 is now a very much abridged version of the previous chapters 5 on the thermal history of the Universe and 6 on particles and symmetries. This reduction was

necessary in order to make space for all the new matter. What has disappeared was standard particle physics which should be known to the particle physics readership (the astronomers do not care).

In Chapter 7 on inflation the currently popular single-field model is rewritten and much expanded.

In Chapter 8 on the cosmic microwave background the emphasis has been shifted from the WMAP satellite to Planck, and all the parameter values have been updated.

Dark matter has been dedicated its own Chapter 9, very much enlarged from the material in the previous edition, and concentrated on the abundant observations of its gravitational effects. Very little is said about dark matter candidates since that is all speculative.

In the very short Chapter 10 on cosmic structures there is nothing new added.

Finally Chapter 11 is devoted to dark energy, almost entirely new in this edition.

I have enjoyed the support of the Magnus Ehrnrooth Foundation for the work on this Edition.

Matts Roos
Helsinki, May 2014

1

From Newton to Hubble

The history of ideas on the structure and origin of the Universe shows that humankind has always put itself at the center of creation. As astronomical evidence has accumulated, these anthropocentric convictions have had to be abandoned one by one. From the natural idea that the solid Earth is at rest and the celestial objects all rotate around us, we have come to understand that we inhabit an average-sized planet orbiting an average-sized sun, that the Solar System is in the periphery of a rotating galaxy of average size, flying at hundreds of kilometres per second towards an unknown goal in an immense Universe, containing billions of similar galaxies.

Cosmology aims to explain the origin and evolution of the entire contents of the Universe, the underlying physical processes, and thereby to obtain a deeper understanding of the laws of physics assumed to hold throughout the Universe. Unfortunately, we have only one universe to study, the one we live in, and we cannot make experiments with it, only observations. This puts serious limits on what we can learn about the origin. If there are other universes we will never know.

Although the history of cosmology is long and fascinating, we shall not trace it in detail, nor any further back than Newton, accounting (in Section 1.1) only for those ideas which have fertilized modern cosmology directly, or which happened to be right although they failed to earn timely recognition. In the early days of cosmology, when little was known about the Universe, the field was really just a branch of philosophy.

Having a rigid Earth to stand on is a very valuable asset. How can we describe motion except in relation to a fixed point? Important understanding has come from the study of inertial systems, in uniform motion with respect to one another. From the work of Einstein on inertial systems, the theory of special relativity was born. In Section 1.2 we discuss inertial frames, and see how expansion and contraction are natural consequences of the homogeneity and isotropy of the Universe.

A classic problem is why the night sky is dark and not blazing like the disc of the Sun, as simple theory in the past would have it. In Section 1.3 we shall discuss this so-called Olbers' paradox, and the modern understanding of it.

The beginning of modern cosmology may be fixed at the publication in 1929 of Hubble's law, which was based on observations of the redshift of spectral lines from remote galaxies. This was subsequently interpreted as evidence for the expansion of the Universe, thus ruling out a static Universe and thereby setting the primary requirement on theory. This will be explained in Section 1.4. In Section 1.5 we turn to determinations of cosmic timescales and the implications of Hubble's law for our knowledge of the age of the Universe.

In Section 1.6 we describe Newton's theory of gravitation, which is the earliest explanation of a gravitational force. We shall 'modernize' it by introducing Hubble's law into it. In fact, we shall see that this leads to a cosmology which already contains many features of current Big Bang cosmologies.

1.1 Historical Cosmology

At the time of *Isaac Newton* (1642–1727) the heliocentric Universe of *Nicolaus Copernicus* (1473–1543), *Galileo Galilei* (1564–1642) and *Johannes Kepler* (1571–1630) had been accepted, because no sensible description of the motion of the planets could be found if the Earth was at rest at the center of the Solar System. Humankind was thus dethroned to live on an average-sized planet orbiting around an average-sized sun.

The stars were understood to be suns like ours with fixed positions in a static Universe. The Milky Way had been resolved into an accumulation of faint stars with the telescope of Galileo. The *anthropocentric view* still persisted, however, in locating the Solar System at the center of the Universe.

Newton's Cosmology. The first theory of gravitation appeared when Newton published his *Philosophiae Naturalis Principia Mathematica* in 1687. With this theory he could explain the empirical laws of Kepler: that the planets moved in elliptical orbits with the Sun at one of the focal points. An early success of this theory came when *Edmund Halley* (1656–1742) successfully predicted that the comet sighted in 1456, 1531, 1607 and 1682 would return in 1758. Actually, the first observation confirming the heliocentric theory came in 1727 when *James Bradley* (1693–1762) discovered the aberration of starlight, and explained it as due to the changes in the velocity of the Earth in its annual orbit. In our time, Newton's theory of gravitation still suffices to describe most of planetary and satellite mechanics, and it constitutes the nonrelativistic limit of Einstein's relativistic theory of gravitation.

Newton considered the stars to be suns evenly distributed throughout infinite space in spite of the obvious concentration of stars in the Milky Way. A distribution is called *homogeneous* if it is uniformly distributed, and it is called *isotropic* if it has the same properties in all spatial directions. Thus in a homogeneous and isotropic space the distribution of matter would look the same to observers located anywhere—no point would be preferential. Each local region of an isotropic universe contains

information which remains true also on a global scale. Clearly, matter introduces lumpiness which grossly violates homogeneity on the scale of stars, but on some larger scale isotropy and homogeneity may still be a good approximation. Going one step further, one may postulate what is called the *cosmological principle*, or sometimes the *Copernican principle*.

The Universe is homogeneous and isotropic in three-dimensional space, has always been so, and will always remain so.

It has always been debated whether this principle is true, and on what scale. On the galactic scale visible matter is lumpy, and on larger scales galaxies form gravitationally bound clusters and narrow strings separated by voids. But galaxies also appear to form loose groups of three to five or more galaxies. Several surveys have now reached agreement that the distribution of these galaxy groups appears to be homogeneous and isotropic within a sphere of 170 Mpc radius [1]. This is an order of magnitude larger than the supercluster to which our Galaxy and our local galaxy group or Local Supercluster (LSC) belong, and which is centered in the constellation of Virgo. Based on his theory of gravitation, Newton formulated a cosmology in 1691. Since all massive bodies attract each other, a finite system of stars distributed over a finite region of space should collapse under their mutual attraction. But this was not observed, in fact the stars were known to have had fixed positions since antiquity, and Newton sought a reason for this stability. He concluded, erroneously, that the self-gravitation within a finite system of stars would be compensated for by the attraction of a sufficient number of stars outside the system, distributed evenly throughout infinite space. However, the total number of stars could not be infinite because then their attraction would also be infinite, making the static Universe unstable. It was understood only much later that the addition of external layers of stars would have no influence on the dynamics of the interior. The right conclusion is that the Universe cannot be static, an idea which would have been too revolutionary at the time.

Newton's contemporary and competitor *Gottfried Wilhelm von Leibnitz* (1646–1716) also regarded the Universe to be spanned by an abstract infinite space, but in contrast to Newton he maintained that the stars must be infinite in number and distributed all over space, otherwise the Universe would be bounded and have a center, contrary to contemporary philosophy. Finiteness was considered equivalent to boundedness, and infinity to unboundedness.

Rotating Galaxies. The first description of the Milky Way as a rotating galaxy can be traced to *Thomas Wright* (1711–1786), who wrote *An Original Theory or New Hypothesis of the Universe* in 1750, suggesting that the stars are

all moving the same way and not much deviating from the same plane, as the planets in their heliocentric motion do round the solar body.

Wright's galactic picture had a direct impact on *Immanuel Kant* (1724–1804). In 1755 Kant went a step further, suggesting that the diffuse nebulae which Galileo had already observed could be distant galaxies rather than nearby clouds of

incandescent gas. This implied that the Universe could be homogeneous on the scale of galactic distances in support of the cosmological principle.

Kant also pondered over the reason for transversal velocities such as the movement of the Moon. If the Milky Way was the outcome of a gaseous nebula contracting under Newton's law of gravitation, why was all movement not directed towards a common center? Perhaps there also existed repulsive forces of gravitation which would scatter bodies onto trajectories other than radial ones, and perhaps such forces at large distances would compensate for the infinite attraction of an infinite number of stars? Note that the idea of a contracting gaseous nebula constituted the first example of a nonstatic system of stars, but at galactic scale with the Universe still static.

Kant thought that he had settled the argument between Newton and Leibnitz about the finiteness or infiniteness of the system of stars. He claimed that either type of system embedded in an infinite space could not be stable and homogeneous, and thus the question of infinity was irrelevant. Similar thoughts can be traced to the scholar *Yang Shen* in China at about the same time, then unknown to Western civilization [2].

The infinity argument was, however, not properly understood until *Bernhard Riemann* (1826–1866) pointed out that the world could be *finite* yet *unbounded*, provided the geometry of the space had a positive curvature, however small. On the basis of Riemann's geometry, *Albert Einstein* (1879–1955) subsequently established the connection between the geometry of space and the distribution of matter.

Kant's repulsive force would have produced trajectories in random directions, but all the planets and satellites in the Solar System exhibit transversal motion in one and the same direction. This was noticed by *Pierre Simon de Laplace* (1749–1827), who refuted Kant's hypothesis by a simple probabilistic argument in 1825: the observed movements were just too improbable if they were due to random scattering by a repulsive force. Laplace also showed that the large transversal velocities and their direction had their origin in the rotation of the primordial gaseous nebula and the law of conservation of angular momentum. Thus no repulsive force is needed to explain the transversal motion of the planets and their moons, no nebula could contract to a point, and the Moon would not be expected to fall down upon us.

This leads to the question of the origin of time: what was the first cause of the rotation of the nebula and when did it all start? This is the question modern cosmology attempts to answer by tracing the evolution of the Universe backwards in time and by reintroducing the idea of a repulsive force in the form of a cosmological constant needed for other purposes.

Black Holes. The implications of Newton's gravity were quite well understood by *John Michell* (1724–1793), who pointed out in 1783 that a sufficiently massive and compact star would have such a strong gravitational field that nothing could escape from its surface. Combining the corpuscular theory of light with Newton's theory, he found that a star with the solar density and escape velocity c would have a radius of $486R_{\odot}$ and a mass of 120 million solar masses. This was the first mention of a type of star much later to be called a *black hole* (to be discussed in Section 3.4). In 1796 Laplace independently presented the same idea.

Galactic and Extragalactic Astronomy. Newton should also be credited with the invention of the reflecting telescope—he even built one—but the first one of importance was built one century later by *William Herschel* (1738–1822). With this instrument, observational astronomy took a big leap forward: Herschel and his son John could map the nearby stars well enough in 1785 to conclude correctly that the Milky Way was a disc-shaped star system. They also concluded erroneously that the Solar System was at its center, but many more observations were needed before it was corrected. Herschel made many important discoveries, among them the planet Uranus, and some 700 binary stars whose movements confirmed the validity of Newton’s theory of gravitation outside the Solar System. He also observed some 250 diffuse nebulae, which he first believed were distant galaxies, but which he and many other astronomers later considered to be nearby incandescent gaseous clouds belonging to our Galaxy. The main problem was then to explain why they avoided the directions of the galactic disc, since they were evenly distributed in all other directions.

The view of Kant that the nebulae were distant galaxies was also defended by *Johann Heinrich Lambert* (1728–1777). He came to the conclusion that the Solar System along, with the other stars in our Galaxy, orbited around the galactic center, thus departing from the heliocentric view. The correct reason for the absence of nebulae in the galactic plane was only given by *Richard Anthony Proctor* (1837–1888), who proposed the presence of interstellar dust. The arguments for or against the interpretation of nebulae as distant galaxies nevertheless raged throughout the 19th century because it was not understood how stars in galaxies more luminous than the whole galaxy could exist—these were observations of supernovae. Only in 1925 did *Edwin P. Hubble* (1889–1953) resolve the conflict indisputably by discovering Cepheids and ordinary stars in nebulae, and by determining the distance to several galaxies, among them the celebrated M31 galaxy in the *Andromeda*. Although this distance was off by a factor of two, the conclusion was qualitatively correct.

In spite of the work of Kant and Lambert, the heliocentric picture of the Galaxy—or almost heliocentric since the Sun was located quite close to Herschel’s galactic center—remained long into our century. A decisive change came with the observations in 1915–1919 by *Harlow Shapley* (1895–1972) of the distribution of *globular clusters* hosting 10^5 – 10^7 stars. He found that perpendicular to the galactic plane they were uniformly distributed, but along the plane these clusters had a distribution which peaked in the direction of the Sagittarius. This defined the center of the Galaxy to be quite far from the Solar System: we are at a distance of about two-thirds of the galactic radius. Thus the anthropocentric world picture received its second blow—and not the last one—if we count Copernicus’s heliocentric picture as the first one. Note that Shapley still believed our Galaxy to be at the center of the astronomical Universe.

The End of Newtonian Cosmology. In 1883 *Ernst Mach* (1838–1916) published a historical and critical analysis of mechanics in which he rejected Newton’s concept of an absolute space, precisely because it was unobservable. Mach demanded that the laws of physics should be based only on concepts which could be related to observations. Since motion still had to be referred to some frame at rest, he proposed replacing absolute space by an idealized rigid frame of fixed stars. Thus ‘uniform

motion' was to be understood as motion relative to the whole Universe. Although Mach clearly realized that all motion is relative, it was left to Einstein to take the full step of studying the laws of physics as seen by observers in inertial frames in relative motion with respect to each other.

Einstein published his General Theory of Relativity in 1917, but the only solution he found to the highly nonlinear differential equations was that of a static Universe. This was not so unsatisfactory though, because the then known Universe comprised only the stars in our Galaxy, which indeed was seen as static, and some nebulae of ill-known distance and controversial nature. Einstein firmly believed in a static Universe until he met Hubble in 1929 and was overwhelmed by the evidence for what was to be called Hubble's law.

Immediately after general relativity became known, *Willem de Sitter* (1872–1934) published (in 1917) another solution, for the case of empty space-time in an exponential state of expansion. In 1922 the Russian meteorologist *Alexandr Friedmann* (1888–1925) found a range of intermediate solutions to the Einstein equation which describe the standard cosmology today. Curiously, this work was ignored for a decade although it was published in widely read journals.

In 1924 Hubble had measured the distances to nine spiral galaxies, and he found that they were extremely far away. The nearest one, M31 in the Andromeda, is now known to be at a distance of 20 galactic diameters (Hubble's value was about 8) and the farther ones at hundreds of galactic diameters. These observations established that the spiral nebulae are, as Kant had conjectured, stellar systems comparable in mass and size with the Milky Way, and their spatial distribution confirmed the expectations of the cosmological principle on the scale of galactic distances.

In 1926–1927 *Bertil Lindblad* (1895–1965) and *Jan Hendrik Oort* (1900–1992) verified Laplace's hypothesis that the Galaxy indeed rotated, and they determined the period to be 10^8 yr and the mass to be about $10^{11} M_{\odot}$. The conclusive demonstration that the Milky Way is an average-sized galaxy, in no way exceptional or central, was given only in 1952 by Walter Baade. This we may count as the third breakdown of the anthropocentric world picture.

The later history of cosmology up until 1990 has been excellently summarized by Peebles [3].

To give the reader an idea of where in the Universe we are, what is nearby and what is far away, some cosmic distances are listed in Table A.1 in the appendix. On a cosmological scale we are not really interested in objects smaller than a galaxy! We generally measure cosmic distances in *parsec* (pc) units (kpc for 10^3 pc and Mpc for 10^6 pc). A parsec is the distance at which one second of arc is subtended by a length equalling the mean distance between the Sun and the Earth. The parsec unit is given in Table A.2 in the appendix, where the values of some useful cosmological and astrophysical constants are listed.

1.2 Inertial Frames and the Cosmological Principle

Newton's first law—the law of inertia—states that a system on which no forces act is either at rest or in uniform motion. Such systems are called *inertial frames*.

Accelerated or rotating frames are not inertial frames. Newton considered that ‘at rest’ and ‘in motion’ implicitly referred to an *absolute space* which was unobservable but which had a real existence independent of humankind. Mach rejected the notion of an empty, unobservable space, and only Einstein was able to clarify the physics of motion of observers in inertial frames.

It may be interesting to follow a nonrelativistic argument about the static or nonstatic nature of the Universe which is a direct consequence of the cosmological principle.

Consider an observer ‘A’ in an inertial frame who measures the density of galaxies and their velocities in the space around him. Because the distribution of galaxies is observed to be homogeneous and isotropic on very large scales (strictly speaking, this is actually true for galaxy groups [1]), he would see the same mean density of galaxies (at one time t) in two different directions \mathbf{r} and \mathbf{r}' :

$$\rho_A(\mathbf{r}, t) = \rho_A(\mathbf{r}', t).$$

Another observer ‘B’ in another inertial frame (see Figure 1.1) looking in the direction \mathbf{r} from her location would also see the same mean density of galaxies:

$$\rho_B(\mathbf{r}', t) = \rho_A(\mathbf{r}, t).$$

The velocity distributions of galaxies would also look the same to both observers, in fact in all directions, for instance in the \mathbf{r}' direction:

$$\mathbf{v}_B(\mathbf{r}', t) = \mathbf{v}_A(\mathbf{r}', t).$$

Suppose that the B frame has the relative velocity $\mathbf{v}_A(\mathbf{r}'', t)$ as seen from the A frame along the radius vector $\mathbf{r}'' = \mathbf{r} - \mathbf{r}'$. If all velocities are nonrelativistic, i.e. small compared with the speed of light, we can write

$$\mathbf{v}_A(\mathbf{r}', t) = \mathbf{v}_A(\mathbf{r} - \mathbf{r}'', t) = \mathbf{v}_A(\mathbf{r}, t) - \mathbf{v}_A(\mathbf{r}'', t).$$

This equation is true only if $\mathbf{v}_A(\mathbf{r}, t)$ has a specific form: it must be proportional to \mathbf{r} ,

$$\mathbf{v}_A(\mathbf{r}, t) = f(t)\mathbf{r}, \tag{1.1}$$

where $f(t)$ is an arbitrary function. Why is this so?

Let this universe start to expand. From the vantage point of A (or B equally well, since all points of observation are equal), nearby galaxies will appear to recede slowly.

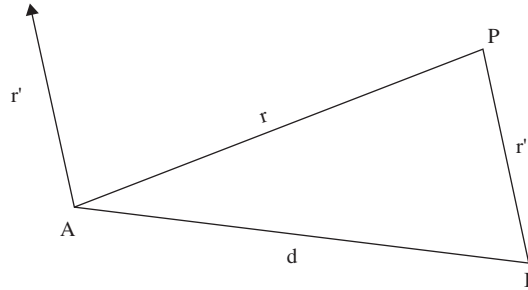


Figure 1.1 Two observers at A and B making observations in the directions \mathbf{r} , \mathbf{r}' .

But in order to preserve uniformity, distant ones must recede faster, in fact their recession velocities must increase linearly with distance. That is the content of Equation (1.1).

If $f(t) > 0$, the Universe would be seen by both observers to expand, each galaxy having a radial velocity proportional to its radial distance r . If $f(t) < 0$, the Universe would be seen to contract with velocities in the reversed direction. Thus we have seen that expansion and contraction are natural consequences of the cosmological principle. If $f(t)$ is a positive constant, Equation (1.1) is Hubble's law.

Actually, it is somewhat misleading to say that the galaxies recede when, rather, it is space itself which expands or contracts. This distinction is important when we come to general relativity.

A useful lesson may be learned from studying the limited gravitational system consisting of the Earth and rockets launched into space. This system is not quite like the previous example because it is not homogeneous, and because the motion of a rocket or a satellite in Earth's gravitational field is different from the motion of galaxies in the gravitational field of the Universe. Thus to simplify the case we only consider radial velocities, and we ignore Earth's rotation. Suppose the rockets have initial velocities low enough to make them fall back onto Earth. The rocket–Earth gravitational system is then *closed* and contracting, corresponding to $f(t) < 0$.

When the kinetic energy is large enough to balance gravity, our idealized rocket becomes a satellite, staying above Earth at a fixed height (real satellites circulate in stable Keplerian orbits at various altitudes if their launch velocities are in the range 8–11 km s⁻¹). This corresponds to the static solution $f(t) = 0$ for the rocket–Earth gravitational system.

If the launch velocities are increased beyond about 11 km s⁻¹, the potential energy of Earth's gravitational field no longer suffices to keep the rockets bound to Earth. Beyond this speed, called the *second cosmic velocity* by rocket engineers, the rockets escape for good. This is an expanding or *open* gravitational system, corresponding to $f(t) > 0$.

The static case is different if we consider the Universe as a whole. According to the cosmological principle, no point is preferred, and therefore there exists no center around which bodies can gravitate in steady-state orbits. Thus the Universe is either expanding or contracting, the static solution being unstable and therefore unlikely.

1.3 Olbers' Paradox

Let us turn to an early problem still discussed today, which is associated with the name of *Wilhelm Olbers* (1758–1840), although it seems to have been known already to Kepler in the 17th century, and a treatise on it was published by *Jean-Philippe Loys de Chéseaux* in 1744, as related in the book by E. Harrison [4]. Why is the night sky dark if the Universe is infinite, static and uniformly filled with stars? They should fill up the total field of visibility so that the night sky would be as bright as the Sun, and we would find ourselves in the middle of a heat bath of the temperature of the surface

of the Sun. Obviously, at least one of the above assumptions about the Universe must be wrong.

The question of the total number of shining stars was already pondered by Newton and Leibnitz. Let us follow in some detail the argument published by Olbers in 1823. The *absolute luminosity* of a star is defined as the amount of luminous energy radiated per unit time, and the *surface brightness* B as luminosity per unit surface. Let the *apparent luminosity* of a star of absolute luminosity L at distance r from an observer be $l = L/4\pi r^2$.

Suppose that the number of stars with average luminosity L is N and their average density in a volume V is $n = N/V$. If the surface area of an average star is A , then its brightness is $B = L/A$. The Sun may be taken to be such an average star, mainly because we know it so well.

The number of stars in a spherical shell of radius r and thickness dr is then $4\pi r^2 n dr$. Their total radiation as observed at the origin of a static universe of infinite extent is then found by integrating the spherical shells from 0 to ∞ :

$$\int_0^{\infty} 4\pi r^2 n l dr = \int_0^{\infty} n L dr = \infty. \quad (1.2)$$

On the other hand, a finite number of visible stars each taking up an angle A/r^2 could cover an infinite number of more distant stars, so it is not correct to integrate r to ∞ . Let us integrate only up to such a distance R that the whole sky of angle 4π would be evenly tiled by the star discs. The condition for this is

$$\int_0^R 4\pi r^2 n \frac{A}{r^2} dr = 4\pi.$$

It then follows that the distance is $R = 1/An$. The integrated brightness from these visible stars alone is then

$$\int_0^R n L dr = L/A, \quad (1.3)$$

or equal to the brightness of the Sun. But the night sky is indeed dark, so we are faced with a paradox.

Olbers' own explanation was that invisible interstellar dust absorbed the light. That would make the intensity of starlight decrease exponentially with distance. But one can show that the amount of dust needed would be so great that the Sun would also be obscured. Moreover, the radiation would heat the dust so that it would start to glow soon enough, thereby becoming visible in the infrared.

A large number of different solutions to this paradox have been proposed in the past, some of the wrong ones lingering on into the present day. Let us here follow a valid line of reasoning due to Lord Kelvin (1824–1907), as retold and improved in a popular book by E. Harrison [4].

A star at distance r covers the fraction $A/4\pi r^2$ of the sky. Multiplying this by the number of stars in the shell, $4\pi r^2 n dr$, we obtain the fraction of the whole sky covered by stars viewed by an observer at the center, $An dr$. Since n is the star count per

volume element, An has the dimensions of number of stars per linear distance. The inverse of this,

$$\ell = 1/An, \quad (1.4)$$

is the mean radial distance between stars, or the *mean free path* of photons emitted from one star and being absorbed in collisions with another. We can also define a mean collision time:

$$\bar{\tau} = \ell/c. \quad (1.5)$$

The value of $\bar{\tau}$ can be roughly estimated from the properties of the Sun, with radius R_\odot and density ρ_\odot . Let the present mean density of luminous matter in the Universe be ρ_0 and the distance to the farthest visible star r_* . Then the collision time inside this volume of size $\frac{4}{3}\pi r_*^3$ is

$$\bar{\tau} \simeq \bar{\tau}_\odot = \frac{1}{A_\odot n c} = \frac{1}{\pi R_\odot^2} \frac{4\pi r_*^3}{3Nc} = \frac{4\rho_\odot R_\odot}{3\rho_0 c}. \quad (1.6)$$

Taking the solar parameters from Table A.2 in the appendix we obtain approximately 10^{23} yr.

The probability that a photon does not collide but arrives safely to be observed by us after a flight distance r can be derived from the assumption that the photon encounters obstacles randomly, that the collisions occur independently and at a constant rate ℓ^{-1} per unit distance. The probability $P(r)$ that the distance to the first collision is r is then given by the exponential distribution

$$P(r) = \ell^{-1} e^{-r/\ell}. \quad (1.7)$$

Thus flight distances much longer than ℓ are improbable.

Applying this to photons emitted in a spherical shell of thickness dr , and integrating the spherical shell from zero radius to r_* , the fraction of all photons emitted in the direction of the center of the sphere and arriving there to be detected is

$$f(r_*) = \int_0^{r_*} \ell^{-1} e^{-r/\ell} dr = 1 - e^{-r_*/\ell}. \quad (1.8)$$

Obviously, this fraction approaches 1 only in the limit of an infinite universe. In that case every point on the sky would be seen to be emitting photons, and the sky would indeed be as bright as the Sun at night. But since this is not the case, we must conclude that r_*/ℓ is small. Thus the reason why the whole field of vision is not filled with stars is that the volume of the presently observable Universe is not infinite, it is in fact too small to contain sufficiently many visible stars.

Lord Kelvin's original result follows in the limit of small r_*/ℓ , in which case

$$f(r_*) \approx r/\ell.$$

The exponential effect in Equation (1.8) was neglected by Lord Kelvin.

We can also replace the mean free path in Equation (1.8) with the collision time [Equation (1.5)], and the distance r_* with the age of the Universe t_0 , to obtain the fraction

$$f(r_*) = g(t_0) = 1 - e^{-t_0/\bar{\tau}}. \quad (1.9)$$

If u_{\odot} is the average radiation density at the surface of the stars, then the radiation density u_0 measured by us is correspondingly reduced by the fraction $g(t_0)$:

$$u_0 = u_{\odot}(1 - e^{-t_0/\bar{\tau}}). \quad (1.10)$$

In order to be able to observe a luminous night sky we must have $u_0 \approx u_{\odot}$, or the Universe must have an age of the order of the collision time, $t_0 \approx 10^{23}$ yr. However, this exceeds all estimates of the age of the Universe by 13 orders of magnitude! Thus the existing stars have not had time to radiate long enough.

What Olbers and many after him did not take into account is that even if the age of the Universe was infinite, the stars do have a finite age and they burn their fuel at well-understood rates.

If we replace 'stars' by 'galaxies' in the above argument, the problem changes quantitatively but not qualitatively. The intergalactic space is filled with radiation from the galaxies, but there is less of it than one would expect for an infinite Universe, at all wavelengths. There is still a problem to be solved, but it is not quite as paradoxical as in Olbers' case.

One explanation is the one we have already met: each star radiates only for a finite time, and each galaxy has existed only for a finite time, whether the age of the Universe is infinite or not. Thus when the time perspective grows, an increasing number of stars become visible because their light has had time to reach us, but at the same time stars which have burned their fuel disappear.

Another possible explanation evokes expansion and special relativity. If the Universe expands, starlight redshifts, so that each arriving photon carries less energy than when it was emitted. At the same time, the volume of the Universe grows, and thus the energy density decreases. The observation of the low level of radiation in the intergalactic space has in fact been evoked as a proof of the expansion.

Since both explanations certainly contribute, it is necessary to carry out detailed quantitative calculations to establish which of them is more important. Most of the existing literature on the subject supports the relativistic effect, but Harrison has shown (and P. S. Wesson [5] has further emphasized) that this is false: the finite lifetime of the stars and galaxies is the dominating effect. The relativistic effect is quantitatively so unimportant that one cannot use it to prove that the Universe is either expanding or contracting.

1.4 Hubble's Law

In the 1920s Hubble measured the spectra of 18 spiral galaxies with a reasonably well-known distance. For each galaxy he could identify a known pattern of atomic spectral lines (from their relative intensities and spacings) which all exhibited a common redward frequency shift by a factor $1 + z$. Using the relation in Equation (1.1) following from the assumption of homogeneity alone,

$$v = cz, \quad (1.11)$$

he could then obtain their velocities with reasonable precision.

The Expanding Universe. The expectation for a stationary universe was that galaxies would be found to be moving about randomly. However, some observations had already shown that most galaxies were redshifted, thus receding, although some of the nearby ones exhibited blueshift. For instance, the nearby Andromeda nebula M31 is approaching us, as its blueshift testifies. Hubble's fundamental discovery was that the velocities of the distant galaxies he had studied increased linearly with distance:

$$v = H_0 r. \quad (1.12)$$

This is called *Hubble's law* and H_0 is called the *Hubble parameter*. For the relatively nearby spiral galaxies he studied, he could only determine the linear, first-order approximation to this function. Although the linearity of this law has been verified since then by the observations of hundreds of galaxies, it is not excluded that the true function has terms of higher order in r . Later on we shall introduce a second-order correction.

The message of Hubble's law is that the Universe is expanding, and this general expansion is called the *Hubble flow*. At a scale of tens or hundreds of Mpc the distances to all astronomical objects are increasing regardless of the position of our observation point. It is true that we observe that the galaxies are receding *from us* as if we were at the center of the Universe. However, we learned from studying a homogeneous and isotropic Universe in Figure 1.1 that if observer A sees the Universe expanding with the factor $f(t)$ in Equation (1.1), any other observer B will also see it expanding with the same factor, and the triangle ABP in Figure 1.1 will preserve its form. Thus, taking the cosmological principle to be valid, every observer will have the impression that all astronomical objects are receding from him/her. A homogeneous and isotropic Universe does not have a center. Consequently, we shall usually talk about *expansion velocities* rather than *recession velocities*.

It is surprising that neither Newton nor later scientists, pondering about why the Universe avoided a gravitational collapse, came to realize the correct solution. An expanding universe would be slowed down by gravity, so the inevitable collapse would be postponed until later. It was probably the notion of an infinite scale of time, inherent in a stationary model, which blocked the way to the right conclusion.

Hubble Time and Radius. From Equations (1.11) and (1.12) one sees that the Hubble parameter has the dimension of inverse time. Thus a characteristic timescale for the expansion of the Universe is the *Hubble time*:

$$\tau_H \equiv H_0^{-1} = 9.7778h^{-1} \times 10^9 \text{ yr.} \quad (1.13)$$

Here h is the commonly used dimensionless quantity

$$h = H_0 / (100 \text{ km s}^{-1} \text{ Mpc}^{-1}).$$

The Hubble parameter also determines the size scale of the observable Universe. In time τ_H , radiation travelling with the speed of light c has reached the *Hubble radius*:

$$r_H \equiv \tau_H c = 3000h^{-1} \text{ Mpc.} \quad (1.14)$$

Or, to put it a different way, according to Hubble's nonrelativistic law, objects at this distance would be expected to attain the speed of light, which is an absolute limit in the theory of special relativity.

Combining Equation (1.12) with Equation (1.11), one obtains

$$z = H_0 \frac{r}{c}. \quad (1.15)$$

In the section on Special Relativity we will see limitations to this formula when v approaches c . The redshift z is in fact infinite for objects at distance r_H receding with the speed of light and thus physically meaningless. Therefore no information can reach us from farther away, all radiation is redshifted to infinite wavelengths, and no particle emitted within the Universe can exceed this distance.

The Cosmic Scale. The size of the Universe is unknown and unmeasurable, but if it undergoes expansion or contraction it is convenient to express distances at different epochs in terms of a *cosmic scale* $R(t)$, and denote its present value $R_0 \equiv R(t_0)$. The value of $R(t)$ can be chosen arbitrarily, so it is often more convenient to normalize it to its present value, and thereby define a dimensionless quantity, the *cosmic scale factor*:

$$a(t) \equiv R(t)/R_0. \quad (1.16)$$

The cosmic scale factor affects all distances: for instance the wavelength λ of light emitted at one time t and observed as λ_0 at another time t_0 :

$$\frac{\lambda_0}{R_0} = \frac{\lambda}{R(t)}. \quad (1.17)$$

Let us find an approximation for $a(t)$ at times $t < t_0$ by expanding it to first-order time differences,

$$a(t) \approx 1 - \dot{a}_0(t_0 - t), \quad (1.18)$$

using the notation \dot{a}_0 for $\dot{a}(t_0)$, and $r = c(t_0 - t)$ for the distance to the source. The *cosmological redshift* can be approximated by

$$z = \frac{\lambda_0}{\lambda} - 1 = a^{-1} - 1 \approx \dot{a}_0 \frac{r}{c}. \quad (1.19)$$

Thus $1/1 + z$ is a measure of the scale factor $a(t)$ at the time when a source emitted the now-redshifted radiation. Identifying the expressions for z in Equations (1.18) and (1.15) we find the important relation

$$\dot{a}_0 = \frac{\dot{R}_0}{R_0} = H_0. \quad (1.20)$$

The Hubble Constant. The value of this constant initially found by Hubble was $H_0 = 550 \text{ km s}^{-1} \text{ Mpc}^{-1}$: an order of magnitude too large because his distance measurements were badly wrong. To establish the linear law and to determine the global value of H_0 one needs to be able to measure distances and expansion velocities well

and far out. Distances are precisely measured only to nearby stars which participate in the general rotation of the Galaxy, and which therefore do not tell us anything about cosmological expansion. Even at distances of several Mpc the expansion-independent, transversal *peculiar velocities* of galaxies are of the same magnitude as the Hubble flow. The measured expansion at the Virgo supercluster, 17 Mpc away, is about 1100 km s^{-1} , whereas the peculiar velocities attain 600 km s^{-1} . At much larger distances where the peculiar velocities do not contribute appreciably to the total velocity, for instance at the Coma cluster 100 Mpc away, the expansion velocity is 6900 km s^{-1} and the Hubble flow can be measured quite reliably, but the imprecision in distance measurements becomes the problem. Every procedure is sensitive to small, subtle corrections and to systematic biases unless great care is taken in the reduction and analysis of data.

Notable contributions to our knowledge of H_0 come from *supernovae* observations with the Hubble Space Telescope (HST) [6, 7], from the measurements of the relic *cosmic microwave background* (CMB) radiation temperature and polarization by the (CMB) radiation temperature Planck satellite [9]. Also the observations WMAP9 [8] and the Baryonic Acoustic Oscillations (BAO) in the distribution of galaxies are important, but the values are reported combined with CMB.

The average of all these experiments [6, 8, 9] is

$$h \equiv H_0 / (100 \text{ km s}^{-1} \text{ Mpc}^{-1}) = 0.696 \pm 0.007. \quad (1.21)$$

Statistics. Let us take the meaning of the term ‘test’ from the statistical literature, where it is accurately defined [10]. When the hypothesis under test concerns the value of a parameter, the problems of *parameter estimation* and *hypothesis testing* are related; for instance, good techniques for estimation often lead to analogous testing procedures. The two situations lead, however, to different conclusions, and should not be confused. If nothing is known *a priori* about the parameter involved, it is natural to use the data to estimate it. On the other hand, if a theoretical prediction has been made that the parameter should have a certain value, it may be more appropriate to formulate the problem as a test of whether the data are consistent with this value. In either case, the nature of the problem, estimation or test, must be clear from the beginning and consistent to the end. When two or more independent methods of parameter estimation are compared, one can talk about a *consistency test*.

A good example of this reasoning is offered by the discussion of Hubble’s law. Hubble’s empirical discovery tested the *null hypothesis* that the Universe (out to the probed redshifts) expands. The test is a valid proof of the hypothesis for any value of H_0 that differs from zero at a chosen confidence level, CL%. Thus the value of $H_0 = 0.673$ is unimportant for the test, only its precision 0.012 matters.

1.5 The Age of the Universe

One of the conclusions of Olbers’ paradox was that the Universe could not be eternal, it must have an age much less than 10^{23} yr, or else the night sky would be bright. More recent proofs that the Universe indeed grows older and consequently has a finite

lifetime comes from astronomical observations of many types of extragalactic objects at high redshifts and at different wavelengths: radio sources, X-ray sources, quasars, faint blue galaxies. High redshifts correspond to earlier times, and what are observed are clear changes in the populations and the characteristics as one looks toward earlier epochs. Let us therefore turn to determinations of the age of the Universe.

In Equation (1.13) we defined the Hubble time τ_H , and gave a value for it of the order of 10 billion years. However, τ_H is not the same as the age t_0 of the Universe. The latter depends on the dynamics of the Universe, whether it is expanding forever or whether the expansion will turn into a collapse, and these scenarios depend on how much matter there is and what the geometry of the Universe is, all questions we shall come back to later.

All the large experiments [11] now agree with an average of

$$t_0 = 13.73 \text{ Gyr.} \quad (1.22)$$

Cosmochronology by Radioactive Nuclei. There are several independent techniques, *cosmochronometers*, for determining the age of the Universe. At this point we shall only describe determinations via the cosmochronology of long-lived radioactive nuclei, and via stellar modeling of the oldest stellar populations in our Galaxy and in some other galaxies. Note that the very existence of radioactive nuclides indicates that the Universe cannot be infinitely old and static.

Various nuclear processes have been used to date the age of the Galaxy, t_G , for instance the ‘Uranium clock’. Long-lived radioactive isotopes such as ^{232}Th , ^{235}U , ^{238}U and ^{244}Pu have been formed by fast neutrons from supernova explosions, captured in the envelopes of an early generation of stars. With each generation of star formation, burn-out and supernova explosion, the proportion of metals increases. Therefore the metal-poorest stars found in globular clusters are the oldest.

The proportions of heavy isotopes following a supernova explosion are calculable with some degree of confidence. Since then, they have decayed with their different natural half-lives so that their abundances in the Galaxy today have changed. For instance, calculations of the original ratio $K = ^{235}\text{U}/^{238}\text{U}$ give values of about 1.3 with a precision of about 10%, whereas this ratio on Earth at the present time is $K_0 = 0.00723$.

To compute the age of the Galaxy by this method, we also need the decay constants λ of ^{238}U and ^{235}U which are related to their half-lives:

$$\lambda_{238} = \ln 2 / (4.46 \text{ Gyr}), \quad \lambda_{235} = \ln 2 / (0.7038 \text{ Gyr}).$$

The relation between isotope proportions, decay constants, and time t_G is

$$K = K_0 \exp [(\lambda_{238} - \lambda_{235})t_G]. \quad (1.23)$$

Inserting numerical values one finds $t_G \approx 6.2 \text{ Gyr}$. However, the Solar System is only 4.57 Gyr old, so the abundance of ^{232}Th , ^{235}U and ^{238}U on Earth cannot be expected to furnish a very interesting limit to t_G . Rather, one has to turn to the abundances on the oldest stars in the Galaxy.

The globular clusters (GCs) are roughly spherically distributed stellar systems in the spheroid of the Galaxy. During the majority of the life of a star, it converts hydrogen into helium in its core. Thus the most interesting stars for the determination of t_G are those which have exhausted their supply of hydrogen, and which are located in old, metal-poor GCs, and to which the distance can be reliably determined. A recent age determination gives

$$t_{GC} = 14.61 \pm 0.8 \text{ Gyr.}$$

This includes an estimated age for the Universe when the clusters formed.

Of particular interest is the detection of a spectral line of ^{238}U in the extremely metal-poor star CS 31082-001, which is overabundant in heavy elements. Theoretical nucleosynthesis models for the initial abundances predict that the ratios of neighboring stable and unstable elements should be similar in early stars as well as on Earth. Thus one compares the abundances of the radioactive ^{232}Th and ^{238}U with the neighboring stable elements Os and Ir (^{235}U is now useless, because it has already decayed away on the oldest stars). One result is

$$t_* = 13.5 \pm 2.9 \text{ Gyr.} \quad (1.24)$$

Brightest Cluster Galaxies (BCGs). Another cosmochronometer is offered by the study of elliptical galaxies in BCGs at very large distances. It has been found that BCG colors only depend on their star-forming histories, and if one can trust stellar population synthesis models, one has a cosmochronometer. From recent analyses of BCGs the result is

$$t_{BCG} \gtrsim 12 \text{ Gyr.} \quad (1.25)$$

Allowing 0.5–1.0 Gyr from the Big Bang until galaxies form stars and clusters, all the above estimates agree reasonably with the value in Equation (1.21) (This correction was already included in the value from globular clusters.).

There are many more cosmochronometers making use of well-understood stellar populations at various distances which we shall not refer to here, all yielding ages near those quoted. It is of interest to note that in the past, when the dynamics of the Universe was less well known, the calculated age τ_H was smaller than the value in Equation (1.21), and at the same time the age t_* of the oldest stars was much higher than the value in Equation (1.23). Thus this historical conflict between cosmological and observational age estimates has now disappeared.

Later we will derive a general relativistic formula for t_0 which depends on a few measurable dynamical parameters determined in a combination of supernova analyses, cosmic microwave background analyses and a set of other data.

1.6 Matter in the Universe

Since antiquity the objects in the sky were known by the visible light they emit, absorb or reflect. Stars like the sun shine, planets and moons reflect sunlight, planets around distant stars reveal themselves by obscuration, and intergalactic dust by dimming

absorption. However, there are other kinds of matter than these examples, and there is radiation at other wavelengths than visible light.

Baryonic Matter. Stable matter as we know it is composed of atoms, and the nuclei of atoms are composed of protons and neutrons which are called nucleons or baryons. The protons are stable particles, the neutrons in atomic nuclei are also stable because of the strong interactions between nucleons. Free neutrons are not stable, they decay dominantly into a proton, an electron and an antineutrino within about 885 s. There exist many more kinds of baryons, but they are unstable and do not form matter.

Stars form galaxies, galaxies form clusters and clusters form superclusters and other large-scale structures. Stars form in the regions of galaxies that are the hardest to observe with many of the common tools of astronomy—in dense, cool (10–100 K) clouds of molecular gas detected in relatively ordinary faraway galaxies. From this environment only a small fraction of visible light can escape. Once stars form, the pressure of their radiation expels the gas, and they can then be seen clearly at optical wavelengths. The results point to a continuous fuelling of gas into the star-forming guts of assembling galaxies.

The baryonic matter in stars and other collapsed objects is only a small fraction of the total baryonic content of the Universe. Much more baryonic matter exists in the form of interstellar dust, hot molecular gas and neutral gas within galaxies, mainly ^1H and ^4He , and in the form of intergalactic hot gas and hot diffuse ionized gas in the intergalactic medium (IGM). The amount of nonradiating diffuse components can be inferred from the absorption of radiation from a bright background source such as a quasar, a technique which is extremely sensitive. Most of the baryonic matter resides outside bound structures, in galaxy groups and in galactic halos.

Current observations of baryons extend from the present-day Solar System to the earliest and most distant galaxies which formed when their age was only 5% of the Universe's present age. About one-fifth of the large galaxies formed within the Universe's first four billion years; 50% of the galaxies had formed by the time the Universe was seven billion years old.

The electromagnetic radiation that stars emit covers all frequencies, not only as visible light but as infrared light, ultraviolet light, X-rays and gamma rays. The most extreme sources of radiation are the *Gamma Ray Bursts (GRB)* from Active Galactic Nuclei (AGN). The nuclear and atomic processes in stars also produce particle emissions: electrons, positrons, neutrinos, antineutrinos and cosmic rays.

There also exists baryonic antimatter, but not on Earth, and there is very little evidence for its presence elsewhere in the Galaxy. That does not mean that antibaryons are pure fiction: they are readily produced in particle accelerators and in violent astrophysical events. However, in an environment of matter, antibaryons rapidly meet baryons and annihilate each other. The asymmetry in the abundance of matter and antimatter is surprising and needs an explanation. We shall deal with that in a later section.

We shall also later see how the baryons came to be the stable end products of the Big Bang Nucleosynthesis and how the mean baryon density in the Universe today is determined from the same set of data as is the age of the Universe.

Supernovae and Neutron Stars. Occasionally, a very bright *supernova* explosion can be seen in some galaxy. These events are very brief (one month) and very rare: historical records show that in our Galaxy they have occurred only every 300 yr. The most recent nearby supernova occurred in 1987 (code name SN1987A), not exactly in our Galaxy but in our small satellite, the Large Magellanic Cloud (LMC). Since it has now become possible to observe supernovae in very distant galaxies, one does not have to wait 300 yr for the next one.

The physical reason for this type of explosion (a Type SNII supernova) is the accumulation of Fe group elements at the core of a massive red giant star of size $8\text{--}200M_{\odot}$, which has already burned its hydrogen, helium and other light elements.

Another type of explosion (a Type SNIa supernova) occurs in binary star systems, composed of a heavy white dwarf and a red giant star. White dwarfs have masses of the order of the Sun, but sizes of the order of Earth, whereas red giants are very large but contain very little mass. The dwarf then accretes mass from the red giant due to its much stronger gravitational field.

As long as the fusion process in the dwarf continues to burn lighter elements to Fe group elements, first the gas pressure and subsequently the electron degeneracy pressure balance the gravitational attraction. But when a rapidly burning dwarf star reaches a mass of $1.44M_{\odot}$, the so-called *Chandrasekhar mass*, or in the case of a red giant when the iron core reaches that mass, no force is sufficient to oppose the gravitational collapse. The electrons and protons in the core transform into neutrinos and neutrons, respectively, most of the gravitational energy escapes in the form of neutrinos, and the remainder is a *neutron star* which is stabilized against further gravitational collapse by the degeneracy pressure of the neutrons. As further matter falls in, it bounces against the extremely dense neutron star and travels outwards as energetic shock waves. In the collision between the shock waves and the outer mantle, violent nuclear reactions take place and extremely bright light is generated. This is the supernova explosion visible from very far away. The nuclear reactions in the mantle create all the elements; in particular, the elements heavier than Fe, Ni and Cr on Earth have all been created in supernova explosions in the distant past.

The released energy is always the same since the collapse always occurs at the Chandrasekhar mass, thus in particular the peak brightness of Type Ia supernovae can serve as remarkably precise standard candles visible from very far away. (The term *standard candle* is used for any class of astronomical objects whose intrinsic luminosity can be inferred independently of the observed flux.) Additional information is provided by the color, the spectrum and an empirical correlation observed between the timescale of the supernova light curve and the peak luminosity. The usefulness of supernovae of Type Ia as standard candles is that they can be seen out to great distances, $z \approx 1.0$, and that the internal precision of the method is quite high. At greater distances one can still find supernovae, but Hubble's linear law [Equation (1.15)] is no longer valid.

The SNeIa are the brightest and most homogeneous class of supernovae. (The plural of SN is abbreviated SNe.) Type II are fainter, and show a wider variation in luminosity. Thus they are not standard candles, but the time evolution of their expanding atmospheres provides an indirect distance indicator, useful out to some 200 Mpc.

The composition of neutron stars is not known. The density of their cores is a few times that of matter in terrestrial nuclei, but they contain far more neutrons than protons, and they are strongly degenerate, thus we have no similar baryonic matter to study in the laboratories. They could be dominated by *quark matter* or by excited forms of baryons such as hyperons which are unstable particles in terrestrial conditions.

Dark components. Nonbaryonic forms of matter or energy which are invisible in the electromagnetic spectrum are neutrinos, black holes, dark matter and dark energy. These components will be dedicated considerable space in later Chapters.

1.7 Expansion in a Newtonian World

In this Section we shall use Newtonian mechanics to derive a cosmology without recourse to Einstein's theory. Inversely, this formulation can also be derived from Einstein's theory in the limit of weak gravitational fields.

A system of massive bodies in an attractive Newtonian potential contracts rather than expands. The Solar System has contracted to a stable, gravitationally bound configuration from some form of hot gaseous cloud, and the same mechanism is likely to be true for larger systems such as the Milky Way, and perhaps also for clusters of galaxies. On yet larger scales the Universe expands, but this does not contradict Newton's law of gravitation.

The key question in cosmology is whether the Universe as a whole is a gravitationally bound system in which the expansion will be halted one day. We shall next derive a condition for this from Newtonian mechanics.

Newtonian Mechanics. Consider a galaxy of *gravitating mass* m_G located at a radius r from the center of a sphere of mean density ρ and mass $M = 4\pi r^3 \rho/3$. The gravitational potential of the galaxy is

$$U = -GMm_G/r = -\frac{4}{3}\pi Gm_G\rho r^2, \quad (1.26)$$

where G is the *Newtonian constant* expressing the strength of the gravitational interaction. Thus the galaxy falls towards the center of gravitation, acquiring a radial acceleration

$$\ddot{r} = -GM/r^2 = -\frac{4}{3}\pi G\rho r. \quad (1.27)$$

This is *Newton's law of gravitation*, usually written in the form

$$F = -\frac{GMm_G}{r^2}, \quad (1.28)$$

where F (in old-fashioned parlance) is the force exerted by the mass M on the mass m_G . The negative signs in Equations (1.28)–(1.30) express the attractive nature of gravitation: bodies are forced to move in the direction of decreasing r .

In a universe expanding linearly according to Hubble's law [Equation (1.12)], the kinetic energy T of the galaxy receding with velocity v is

$$T = \frac{1}{2}mv^2 = \frac{1}{2}mH_0^2r^2, \quad (1.29)$$

where m is the *inertial mass* of the galaxy. Although there is no theoretical reason for the inertial mass to equal the gravitational mass (we shall come back to this question later), careful tests have verified the equality to a precision better than a few parts in 10^{13} . Let us therefore set $m_G = m$. Thus the total energy is given by

$$E = T + U = \frac{1}{2}mH_0^2r^2 - \frac{4}{3}\pi Gm\rho r^2 = mr^2 \left(\frac{1}{2}H_0^2 - \frac{4}{3}\pi G\rho \right). \quad (1.30)$$

If the mass density ρ of the Universe is large enough, the expansion will halt. The condition for this to occur is $E = 0$, or from Equation (1.32) this *critical density* is

$$\rho_c = \frac{3H_0^2}{8\pi G} = 1.0539 \times 10^{10} h^2 \text{ eV m}^{-3}. \quad (1.31)$$

The value $h = 0.696$ from Equation (1.21) can be inserted here. A universe with density $\rho > \rho_c$ is called *closed*; with density $\rho < \rho_c$ it is called *open*.

Expansion. Note that r and ρ are time dependent: they scale with the expansion. Denoting their present values r_0 and ρ_0 , one has

$$r(t) = r_0 a(t), \quad \rho(t) = \rho_0 a^{-3}(t). \quad (1.32)$$

The acceleration \ddot{r} in Equation (1.27) can then be replaced by the acceleration of the scale

$$\ddot{a} = \ddot{r}/r_0 = -\frac{4}{3}\pi G\rho_0 a^{-2}. \quad (1.33)$$

Let us use the identity

$$\ddot{a} = \frac{1}{2} \frac{d}{da} \dot{a}^2$$

in Equation (1.33) to obtain

$$d\dot{a}^2 = -\frac{8}{3}\pi G\rho_0 \frac{da}{a^2}.$$

This can be integrated from the present time t_0 to an earlier time t with the result

$$\dot{a}^2(t) - \dot{a}^2(t_0) = \frac{8}{3}\pi G\rho_0 (a^{-1} - 1). \quad (1.34)$$

Let us now introduce the dimensionless *density parameter*:

$$\Omega_0 = \frac{\rho_0}{\rho_c} = \frac{8\pi G\rho_0}{3H_0^2}. \quad (1.35)$$

Substituting Ω_0 into Equation (1.34) and making use of the relation in Equation (1.20), $\dot{a}(t_0) = H_0$, we find

$$\dot{a}^2 = H_0^2 (\Omega_0 a^{-1} - \Omega_0 + 1). \quad (1.36)$$

Thus it is clear that the presence of matter influences the dynamics of the Universe. Without matter, $\Omega_0 = 0$, Equation (1.36) just states that the expansion is constant, $\dot{a} = H_0$, and H_0 could well be zero as Einstein thought. During expansion \dot{a} is positive; during contraction it is negative. In both cases the value of \dot{a}^2 is nonnegative, so it must always be true that

$$1 - \Omega_0 + \Omega_0/a \geq 0. \quad (1.37)$$

Models of Cosmological Evolution. Depending on the value of Ω_0 the evolution of the Universe can take three courses.

- (i) $\Omega_0 < 1$, the mass density is undercritical. As the cosmic scale factor $a(t)$ increases for times $t > t_0$ the term Ω_0/a decreases, but the expression (1.37) stays positive always. Thus this case corresponds to an open, ever-expanding universe, as a consequence of the fact that it is expanding now. In Figure 1.2 the expression in Equation (1.37) is plotted against a as the long-dashed curve for the choice $\Omega_0 = 0.5$.
- (ii) $\Omega_0 = 1$, the mass density is critical. As the scale factor $a(t)$ increases for times $t > t_0$ the expression in Equation (1.37) gradually approaches zero, and the expansion halts. However, this only occurs infinitely late, so it also corresponds to an ever-expanding universe. This case is plotted against a as the short-dashed curve in Figure 1.2. Note that cases (i) and (ii) differ by having different asymptotes. Case (ii) is quite realistic because the observational value of Ω_0 is very close to 1, as we shall see later.

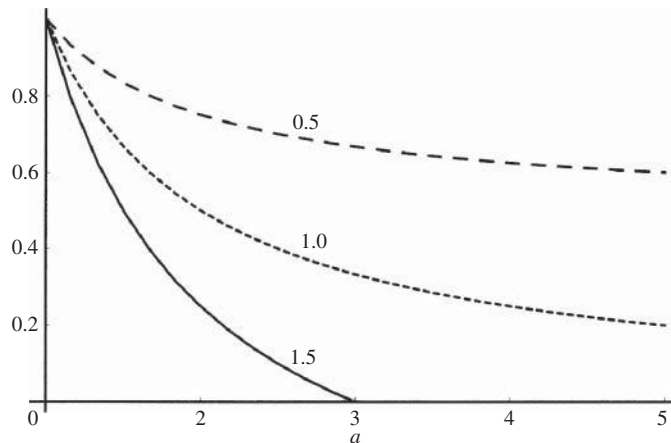


Figure 1.2 Dependence of the expression in Equation (1.37) on the cosmic scale a for an undercritical ($\Omega_0 = 0.5$), critical ($\Omega_0 = 1$) and overcritical ($\Omega_0 = 1.5$) universe. Time starts today at scale $a = 1$ in this picture and increases with a , except for the overcritical case where the Universe arrives at its maximum size, here $a = 3$, whereupon it reverses its direction and starts to shrink.

- (iii) $\Omega_0 > 1$, the mass density is overcritical and the Universe is closed. As the scale factor $a(t)$ increases, it reaches a maximum value a_{mid} when the expression in Equation (1.37) vanishes, and where the rate of increase, \dot{a}_{mid} , also vanishes. But the condition (1.37) must stay true, and therefore the expansion must turn into contraction at a_{mid} . The solid line in Figure 1.2 describes this case for the choice $\Omega_0 = 1.5$, whence $a_{\text{mid}} = 3$. For later times the Universe retraces the solid curve, ultimately reaching scale $a = 1$ again.

This is as far as we can go combining Newtonian mechanics with Hubble's law. We have seen that problems appear when the recession velocities exceed the speed of light, conflicting with special relativity. Another problem is that Newton's law of gravitation knows no delays: the gravitational potential is felt instantaneously over all distances. A third problem with Newtonian mechanics is that the Copernican world, which is assumed to be homogeneous and isotropic, extends up to a finite distance r_0 , but outside that boundary there is nothing. Then the boundary region is characterized by violent inhomogeneity and anisotropy, which are not taken into account. To cope with these problems we must begin to construct a fully relativistic cosmology.

Problems

1. How many revolutions has the Galaxy made since the formation of the Solar System if we take the solar velocity around the galactic center to be 365 km s^{-1} ?
2. Use Equation (1.4) to estimate the mean free path ℓ of photons. What fraction of all photons emitted by stars up to the maximum observed redshift $z = 7$ arrive at Earth?
3. If Hubble had been right that the expansion is given by

$$H_0 = 550 \text{ km s}^{-1} \text{ Mpc}^{-1},$$

how old would the Universe be then [see Equation (1.13)]?

4. What is the present ratio $K_0 = {}^{235}\text{U}/{}^{238}\text{U}$ on a star 10 Gyr old?
5. Prove Newton's theorem that the gravitational force at a radial distance R from the center of a spherical distribution of matter acts as if all the mass inside R were concentrated at a single point at the center. Show also that if the spherical distribution of matter extends beyond R , the force due to the mass outside R vanishes.
6. Estimate the escape velocity from the Galaxy.

References

- [1] Ramella, M., Geller, M. J., Pisani, A. and da Costa, L. N. 2002 *Astron. J.* **123**, 2976.
- [2] Fang Li Zhi and Li Shu Xian 1989 *Creation of the Universe*. World Scientific, Singapore.
- [3] Peebles, P. J. E. 1993 *Principles of physical cosmology*. Princeton University Press, Princeton, NJ.
- [4] Harrison, E. 1987 *Darkness at night*. Harvard University Press, Cambridge, MA.

- [5] Wesson, P. S. 1991 *Astrophys. J.* **367**, 399.
- [6] Efstathiou, G., 2013. Preprint arXiv:1311.3461 [astro-ph.CO].
- [7] Freedman, W. L. *et al.* 2001 *Astrophys. J.* **553**, 47.
- [8] Bennett, C. L. *et al.* 2013 *Astrophys. J. Suppl.* **208**, 20.
- [9] Planck Collaboration: Ade, P. A. R. *et al.* 2014 *Astron. Astrophys* and preprint arXiv:1303.5076 [astro-ph.CO]
- [10] Eadie, W. T., Drijard, D., James, F. E., Roos, M. and Sadoulet, B. 1971 *Statistical methods in experimental physics*. North-Holland, Amsterdam. Second Edition by Frederick James 2006. World Scientific Publishing Co.
- [11] Bennett, C.L. *et al.* 2014, *Astrophys. J.* **794**, 135.

Special Relativity

The foundations of modern cosmology were laid during the second and third decade of the twentieth century: on the theoretical side by Einstein's theory of general relativity, which represented a deep revision of current concepts; and on the observational side by Hubble's discovery of the cosmic expansion, which ruled out a static Universe and set the primary requirement on theory. Space and time are not invariants under Lorentz transformations, their values being different to observers in different inertial frames. Nonrelativistic physics uses these quantities as completely adequate approximations, but in relativistic frame-independent physics we must find invariants to replace them. This chapter begins, in Section 2.1, with Einstein's theory of special relativity, which gives us such invariants.

In Section 2.2 we generalize the metrics in linear spaces to metrics in curved spaces, in particular the Robertson–Walker metric in a four-dimensional manifold. This gives us tools to define invariant distance measures in Section 2.3, which are the key to Hubble's parameter. To conclude we discuss briefly tests of special relativity in Section 2.4.

2.1 Lorentz Transformations

In Einstein's theory of special relativity one studies how signals are exchanged between inertial frames in linear motion with respect to each other with constant velocity. Einstein made two postulates about such frames:

- (i) the results of measurements in different frames must be identical;
- (ii) light travels by a constant speed, c , in vacuo, in all frames.

The first postulate requires that physics be expressed in frame-independent invariants. The latter is actually a statement about the measurement of time in different frames, as we shall see shortly.

Lorentz Transformations. Consider two linear axes x and x' in one-dimensional space, x' being at rest and x moving with constant velocity v in the positive x' direction. Time increments are measured in the two coordinate systems as dt and dt' using two identical clocks. Neither the spatial increments dx and dx' nor the time increments dt and dt' are invariants—they do not obey postulate (i). Let us replace dt and dt' with the temporal distances $c dt$ and $c dt'$ and look for a *linear transformation* between the primed and unprimed coordinate systems, under which the two-dimensional *space-time distance* ds between two *events*,

$$ds^2 = c^2 d\tau^2 = c^2 dt^2 - dx^2 = c^2 dt'^2 - dx'^2, \quad (2.1)$$

is invariant. Invoking the constancy of the speed of light it is easy to show that the transformation must be of the form

$$dx' = \gamma(dx - v dt), \quad c dt' = \gamma(c dt - v dx/c), \quad (2.2)$$

where

$$\gamma = \frac{1}{\sqrt{1 - (v/c)^2}}. \quad (2.3)$$

Equation (2.2) defines the *Lorentz transformation*, after *Hendrik Antoon Lorentz* (1853–1928). Scalar products in this two-dimensional (ct, x) -space are invariants under Lorentz transformations.

Time Dilation. The quantity $d\tau$ in Equation (2.1) is called the *proper time* and ds the *line element*. Note that scalar multiplication in this manifold is here defined in such a way that the products of the spatial components obtain negative signs (sometimes the opposite convention is chosen). (The mathematical term for a many-dimensional space is a *manifold*.)

Since $d\tau^2$ is an invariant, it has the same value in both frames:

$$d\tau'^2 = d\tau^2.$$

While the observer at rest records consecutive ticks on his clock separated by a space-time interval $d\tau = dt'$, she receives clock ticks from the x direction separated by the time interval dt and also by the space interval $dx = v dt$:

$$d\tau = d\tau' = \sqrt{dt^2 - dx^2/c^2} = \sqrt{1 - (v/c)^2} dt. \quad (2.4)$$

In other words, the two inertial coordinate systems are related by a Lorentz transformation

$$dt = \frac{dt'}{\sqrt{1 - (v/c)^2}} \equiv \gamma dt'. \quad (2.5)$$

Obviously, the time interval dt is always longer than the interval dt' , but only noticeably so when v approaches c . This is called the *time dilation effect*.

The time dilation effect has been well confirmed in particle experiments. Muons are heavy, unstable, electron-like particles with well-known lifetimes in the laboratory. However, when they strike Earth with relativistic velocities after having been

produced in cosmic ray collisions in the upper atmosphere, they appear to have a longer lifetime by the factor γ .

Another example is furnished by particles of mass m and charge Q circulating with velocity v in a synchrotron of radius r . In order to balance the centrifugal force the particles have to be subject to an inward-bending magnetic field density B . The classical condition for this is

$$r = mv/QB.$$

The velocity in the circular synchrotron as measured at rest in the laboratory frame is inversely proportional to t , say the time of one revolution. But in the particle rest frame the time of one revolution is shortened to t/γ . When the particle attains relativistic velocities (by traversing accelerating potentials at regular positions in the ring), the magnetic field density B felt by the particle has to be adjusted to match the velocity in the particle frame, thus

$$r = mv\gamma/QB.$$

This equation has often been misunderstood to imply that the mass m increases by the factor γ , whereas only time measurements are affected by γ .

Relativity and Gold. Another example of relativistic effects on the orbits of circulating massive particles is furnished by electrons in Bohr orbits around a heavy nucleus. The effective Bohr radius of an electron is inversely proportional to its mass. Near the nucleus the electrons attain relativistic speeds, the time dilation will cause an apparent increase in the electron mass, more so for inner electrons with larger average speeds. For a 1s shell at the nonrelativistic limit, this average speed is proportional to Z atomic units. For instance, v/c for the 1s electron in Hg is $80/137 = 0.58$, implying a relativistic radial shrinkage of 23%. Because the higher s shells have to be orthogonal against the lower ones, they will suffer a similar contraction. Due to interacting relativistic and shell-structure effects, their contraction can be even larger; for gold, the 6s shell has larger percentage relativistic effects than the 1s shell. The nonrelativistic 5d and 6s orbital energies of gold are similar to the 4d and 5s orbital energies of silver, but the relativistic energies happen to be very different. This is the cause of the chemical difference between silver and gold and also the cause for the distinctive color of gold [2].

Light Cone. The Lorentz transformations [Equations (2.1), (2.2)] can immediately be generalized to three spatial dimensions, where the square of the Pythagorean distance element

$$dl^2 \equiv d\mathbf{l}^2 = dx^2 + dy^2 + dz^2 \quad (2.6)$$

is invariant under rotations and translations in three-space. This is replaced by the four-dimensional space-time of *Hermann Minkowski* (1864–1909), defined by the temporal distance ct and the spatial coordinates x, y, z . An invariant under Lorentz

transformations between frames which are rotated or translated at a constant velocity with respect to each other is then the line element of the *Minkowski metric*

$$ds^2 = c^2 d\tau^2 = c^2 dt^2 - dx^2 - dy^2 - dz^2 = c^2 dt^2 - dl^2. \quad (2.7)$$

The trajectory of a body moving in space-time is called its *world line*. A body at a fixed location in space follows a world line parallel to the time axis and, of course, in the direction of increasing time. A body moving in space follows a world line making a slope with respect to the time axis. Since the speed of a body or a signal travelling from one event to another cannot exceed the speed of light, there is a maximum slope to such world lines. All world lines arriving where we are, here and now, obey this condition. Thus they form a cone in our past, and the envelope of the cone corresponds to signals travelling with the speed of light. This is called the *light cone*.

Two separate events in space-time can be *causally* connected provided their spatial separation $d\mathbf{l}$ and their temporal separation dt (in any frame) obey

$$|d\mathbf{l}/dt| \leq c.$$

Their world line is then inside the light cone. In Figure 2.1 we draw this four-dimensional cone in t, x, y -space (another choice could have been to use the coordinates t, σ, θ). Thus if we locate *our present* event to the apex of the light cone at $t = t_0 = 0$, it can be influenced by world lines from all events inside the *past* light cone for which $ct < 0$, and it can influence all events inside the *future* light cone for which $ct > 0$. Events inside the light cone are said to have *timelike* separation from the present event. Events outside the light cone are said to have *spacelike* separation from the present event: they cannot be causally connected to it. Thus the light cone encloses the *present observable universe*, which consists of all world lines that can in principle be observed. From now on we usually mean the present observable universe when we say simply ‘the Universe’.

For light signals the equality sign above applies so that the proper time interval in Equation (2.7) vanishes

$$d\tau = 0.$$

Events on the light cone are said to have *null* or *lightlike* separation.

Redshift and Scale Factor. The light emitted by stars is caused by atomic transitions with emission spectra containing sharp spectral lines. Similarly, hot radiation traversing cooler matter in stellar atmospheres excites atoms at sharply defined wavelengths, producing characteristic dark absorption lines in the continuous regions of the emission spectrum. The radiation that was emitted by stars and distant galaxies with a wavelength $\lambda_{\text{rest}} = c/\nu_{\text{rest}}$ at time t in their rest frame will have its wavelength stretched by the cosmological expansion to λ_{obs} when observed on Earth. Since the Universe expands, this shift is in the red direction, $\lambda_{\text{obs}} > \lambda_{\text{rest}}$, and it is therefore called a *redshift*, denoted

$$z = \frac{\lambda_{\text{obs}} - \lambda_{\text{rest}}}{\lambda_{\text{rest}}}. \quad (2.8)$$

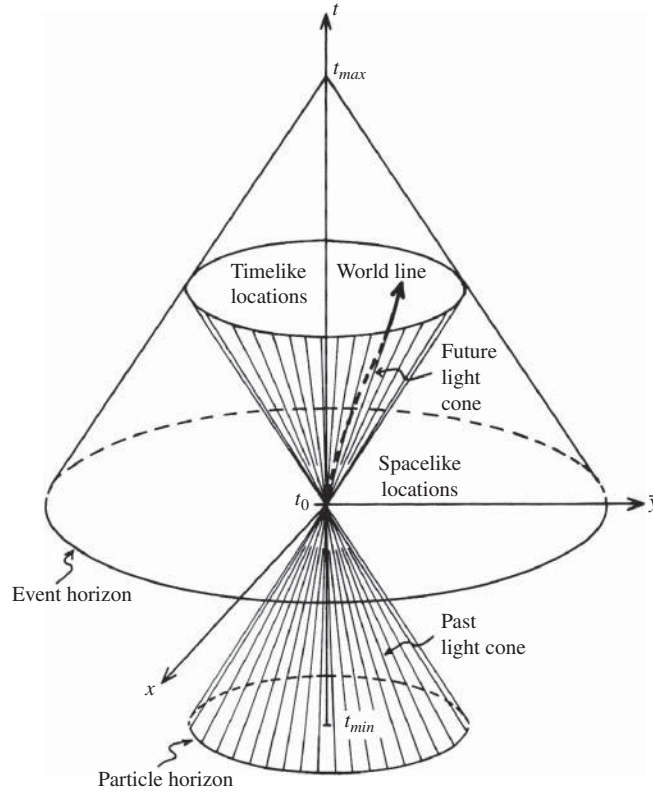


Figure 2.1 Light cone in x, y, t -space. An event which is at the origin $x = y = 0$ at the present time t_0 will follow some world line into the future, always remaining inside the future light cone. All points on the world line are at timelike locations with respect to the spatial origin at t_0 . World lines for light signals emitted from (received at) the origin at t_0 will propagate on the envelope of the future (past) light cone. No signals can be sent to or received from spacelike locations. The space in the past from which signals can be received at the present origin is restricted by the particle horizon at t_{min} , the earliest time under consideration. The event horizon restricts the space which can at present be in causal relation to the present spatial origin at some future time t_{max} .

The redshift is denoted by the letter z , not to be confused with the coordinate z in Equations (2.6) and (2.7).

The ratio of wavelengths actually measured by the terrestrial observer is then

$$1 + z = \frac{\lambda_{obs}}{\lambda_{rest}} = \frac{R_0}{R(t)} = \frac{1}{a(t)}. \quad (2.9)$$

It should be stressed that the cosmological redshift is not caused by the velocities of receding objects, but by the increase in scale $a(t)$ since time t . A kinematic effect can be observed in the spectra of nearby stars and galaxies, for which their peculiar motion is more important than the effect of the cosmological expansion. This may give rise

to a *Doppler redshift* for a receding source, and to a corresponding *blueshift* for an approaching source.

Actually, the light cones in Figure 2.1 need to be modified for an expanding universe. A scale factor $a(t)$ that increases with time implies that light will travel a distance greater than ct during time t . Consequently, the straight lines defining the cone will be curved outwards.

2.2 Metrics of Curved Space-time

In Newton's time the laws of physics were considered to operate in a *flat Euclidean space*, in which spatial distance could be measured on an infinite and immovable three-dimensional grid, and time was a parameter marked out on a linear scale running from infinite past to infinite future. But Newton could not answer the question of how to identify which inertial frame was at rest relative to this absolute space. In his days the solar frame could have been chosen, but today we know that the Solar System orbits the Galactic center, the Galaxy is in motion relative to the local galaxy group, which in turn is in motion relative to the Hydra–Centaurus cluster, and the whole Universe is expanding.

The geometry of curved spaces was studied in the nineteenth century by Gauss, Riemann and others. Riemann realized that Euclidean geometry was just a particular choice suited to flat space, but not necessarily correct in the space we inhabit. And Mach realized that one had to abandon the concept of absolute space altogether. Einstein learned about *tensors* from his friend Marcel Grossman, and used these key quantities to go from flat Euclidean three-dimensional space to curved Minkowskian four-dimensional space in which physical quantities are described by invariants. Tensors are quantities which provide generally valid relations between different four-vectors.

Euclidean Space. Let us consider how to describe distance in three-space. The path followed by a free body obeying Newton's first law of motion can suitably be described by expressing its spatial coordinates as functions of time: $x(t)$, $y(t)$, $z(t)$. Time is then treated as an absolute parameter and not as a coordinate. This path represents the shortest distance between any two points along it, and it is called a *geodesic* of the space. As is well known, in Euclidean space the geodesics are straight lines. Note that the definition of a geodesic does not involve any particular coordinate system.

If one replaces the components x , y , z of the distance vector \mathbf{l} by x^1 , x^2 , x^3 , this permits a more compact notation of the Pythagorean squared distance l^2 in the *metric* [Equation (2.6)]

$$l^2 = (x^1)^2 + (x^2)^2 + (x^3)^2 = \sum_{i,j=1}^3 g_{ij} x^i x^j \equiv g_{ij} x^i x^j. \quad (2.10)$$

The quantities g_{ij} are the nine components of the *metric tensor* \mathbf{g} , which contains all the information about the intrinsic geometry of this three-space. In the last step we

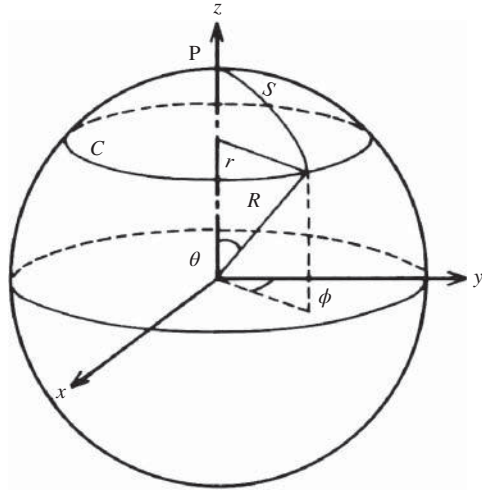


Figure 2.2 A two-sphere on which points are specified by coordinates (θ, ϕ) .

have used the convention to leave out the summation sign; it is then implied that summation is carried out over repeated indices. One commonly uses Roman letters in the indices when only the spatial components x^i , $i = 1, 2, 3$, are implied, and Greek letters when all the four space-time coordinates x^μ , $\mu = 0, 1, 2, 3$, are implied. Orthogonal coordinate systems have diagonal metric tensors and this is all that we will encounter. The components of g in flat Euclidean three-space are

$$g_{ij} = \delta_{ij},$$

where δ_{ij} is the usual Kronecker delta.

The same flat space could equally well be mapped by, for example, spherical or cylindrical coordinates. The components g_{ij} of the metric tensor would be different, but Equation (2.10) would hold unchanged. For instance, choosing spherical coordinates R, θ, ϕ as in Figure 2.2,

$$x = R \sin \theta \sin \phi, \quad y = R \sin \theta \cos \phi, \quad z = R \cos \theta, \quad (2.11)$$

dl^2 takes the explicit form

$$dl^2 = dR^2 + R^2 d\theta^2 + R^2 \sin^2 \theta d\phi^2. \quad (2.12)$$

Geodesics in this space obey Newton's first law of motion, which may be written as

$$\ddot{\mathbf{R}} = 0. \quad (2.13)$$

Dots indicate time derivatives.

Minkowski Space-time. In special relativity, symmetry between spatial coordinates and time is achieved, as is evident from the Minkowski metric (2.7) describing a flat space-time in four Cartesian coordinates. In tensor notation the Minkowski metric

includes the coordinate $dx^0 \equiv c dt$ so that the invariant line element in Equation (2.7) can be written

$$ds^2 = c^2 d\tau^2 = g_{\mu\nu} dx^\mu dx^\nu. \quad (2.14)$$

The 16 components of \mathbf{g} in flat Minkowski space-time are given by the diagonal matrix $\eta_{\mu\nu}$, a generalization of the Kronecker delta function to four-space-time,

$$\eta_{00} = 1, \quad \eta_{jj} = -1, \quad j = 1, 2, 3, \quad (2.15)$$

all nondiagonal components vanishing. The choice of signs in the definition of $\eta_{\mu\nu}$ is not standardized in the literature, but we shall use Equation (2.15).

The path of a body, or its world line, is then described by the four coordinate functions $x(\tau)$, $y(\tau)$, $z(\tau)$, $t(\tau)$, where the proper time τ is a new absolute parameter, an invariant under Lorentz transformations. A geodesic in the Minkowski space-time is also a straight line, given by the equations

$$\frac{d^2 t}{d\tau^2} = 0, \quad \frac{d^2 \mathbf{R}}{d\tau^2} = 0. \quad (2.16)$$

In the spherical coordinates [Equation (2.11)] the Minkowski metric [Equation (2.7)] takes the form

$$ds^2 = c^2 dt^2 - dl^2 = c^2 dt^2 - dR^2 - R^2 d\theta^2 - R^2 \sin^2 \theta d\phi^2. \quad (2.17)$$

An example of a curved space is the two-dimensional surface of a sphere with radius R obeying the equation

$$x^2 + y^2 + z^2 = R^2. \quad (2.18)$$

This surface is called a two-sphere.

Combining Equations (2.6) and (2.18) we see that one coordinate is really superfluous, for instance z , so that the spatial metric [Equation (2.6)] can be written

$$dl^2 = dx^2 + dy^2 + \frac{(x dx + y dy)^2}{R^2 - x^2 - y^2}. \quad (2.19)$$

This metric describes spatial distances on a two-dimensional surface embedded in three-space, but the third dimension is not needed to measure a distance on the surface. Note that R is not a third coordinate, but a constant everywhere on the surface.

Thus measurements of distances depend on the geometric properties of space, as has been known to navigators ever since Earth was understood to be spherical. The geodesics on a sphere are great circles, and the metric is

$$dl^2 = R^2 d\theta^2 + R^2 \sin^2 \theta d\phi^2. \quad (2.20)$$

Near the poles where $\theta = 0^\circ$ or $\theta = 180^\circ$ the local distance would depend very little on changes in longitude ϕ . No point on this surface is preferred, so it can correspond to a Copernican homogeneous and isotropic two-dimensional universe which is unbounded, yet finite.

Let us write Equation (2.20) in the matrix form

$$dl^2 = (d\theta \quad d\phi) \mathbf{g} \begin{pmatrix} d\theta \\ d\phi \end{pmatrix}, \quad (2.21)$$

where the metric matrix is

$$\mathbf{g} = \begin{pmatrix} R^2 & 0 \\ 0 & R^2 \sin^2 \theta \end{pmatrix}. \quad (2.22)$$

The ‘two-volume’ or area A of the two-sphere in Figure 2.2 can then be written

$$A = \int_0^{2\pi} d\phi \int_0^\pi d\theta \sqrt{\det \mathbf{g}} = \int_0^{2\pi} d\phi \int_0^\pi d\theta R^2 \sin \theta = 4\pi R^2, \quad (2.23)$$

as expected.

In Euclidean three-space parallel lines of infinite length never cross, but this could not be proved in Euclidean geometry, so it had to be asserted without proof, the *parallel axiom*. The two-sphere belongs to the class of Riemannian curved spaces which are locally flat: a small portion of the surface can be approximated by its tangential plane. Lines in this plane which are parallel locally do cross when extended far enough, as required for geodesics on the surface of a sphere.

Gaussian Curvature. The deviation of a curved surface from flatness can also be determined from the length of the circumference of a circle. Choose a point ‘P’ on the surface and draw the locus corresponding to a fixed distance s from that point. If the surface is flat, a plane, the locus is a circle and s is its radius. On a two-sphere of radius R the locus is also a circle, see Figure 2.2, but the distance s is measured along a geodesic. The angle subtended by s at the center of the sphere is s/R , so the radius of the circle is $r = R \sin(s/R)$. Its circumference is then

$$C = 2\pi R \sin(s/R) = 2\pi s \left(1 - \frac{s^2}{6R^2} + \dots \right). \quad (2.24)$$

Carl Friedrich Gauss (1777–1855) discovered an invariant characterizing the curvature of two-surfaces, the *Gaussian curvature* K . Although K can be given by a completely general formula independent of the coordinate system (see, e.g., [1]), it is most simply described in an orthogonal system x, y . Let the radius of curvature along the x -axis be $R_x(x)$ and along the y -axis be $R_y(y)$. Then the Gaussian curvature at the point (x_0, y_0) is

$$K = 1/R_x(x_0)R_y(y_0). \quad (2.25)$$

On a two-sphere $R_x = R_y = R$, so $K = R^{-2}$ everywhere. Inserting this into Equation (2.24) we obtain, in the limit of small s ,

$$K = \frac{3}{\pi} \lim_{s \rightarrow 0} \left(\frac{2\pi s - C}{s^3} \right). \quad (2.26)$$

This expression is true for any two-surface, and it is in fact the only invariant that can be defined.

Whether we live in three or more dimensions, and whether our space is flat or curved, is really a physically testable property of space. Gauss actually proceeded to investigate this by measuring the angles in a triangle formed by three distant mountain peaks. If space were Euclidean the value would be 180° , but if the surface had

positive curvature like a two-sphere the angles would add up to more than 180° . Correspondingly, the angles on a saddle surface with negative curvature would add up to less than 180° . This is illustrated in Figures 2.3 and 2.4. The precision in Gauss's time was, however, not good enough to exhibit any disagreement with the Euclidean value.

Comoving Coordinates. If the two-sphere with surface [Equation (2.18)] and constant radius R were a balloon expanding with time we replace R by the expansion scale $a(t)$, defined in Equation (1.16). Points on the surface of the balloon would find their mutual distances scaled by $a(t)$ relative to a time t_0 when the radius was $R_0 = 1$.

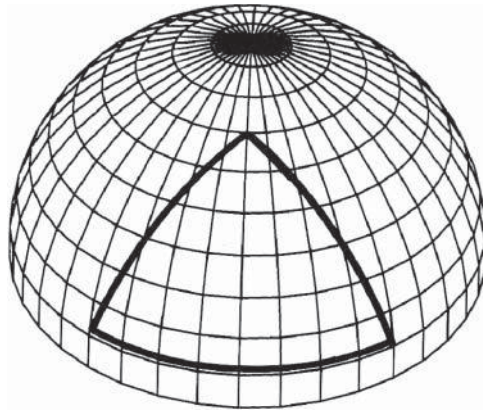


Figure 2.3 The angles in a triangle on a surface with positive curvature add up to more than 180° .

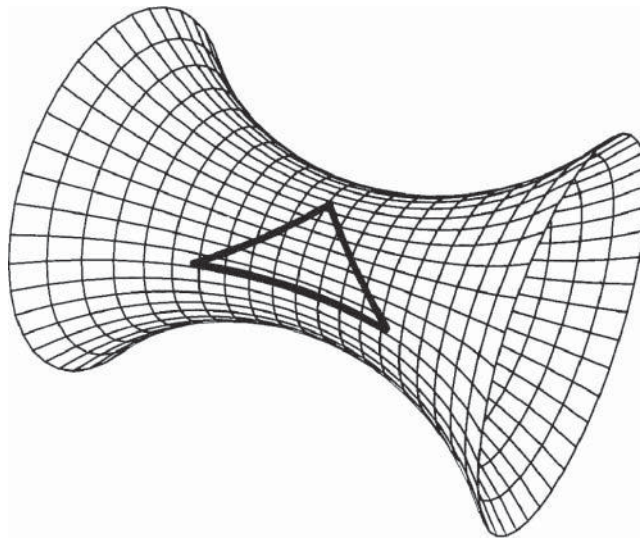


Figure 2.4 The angles in a triangle on a surface with negative curvature add up to less than 180° .

An observer located at any one point would see all the other points recede radially. This is exactly how we see distant galaxies except that we are not on a two-sphere but, as we shall see, on a spatially curved three-surface with the cosmic scale factor $a(t)$.

Suppose this observer wants to make a map of all the points on the expanding surface. It is then no longer convenient to use coordinates dependent on $a(t)$ as in Equations (2.12) and (2.20), because the map would quickly be outdated. Instead it is convenient to factor out the cosmic expansion and replace a by $a(t)\sigma$, where σ is a dimensionless *comoving* coordinate, thus

$$dl^2 = a^2(t)(d\sigma^2 + \sigma^2 d\theta^2 + \sigma^2 \sin^2\theta d\phi^2). \quad (2.27)$$

Returning to the space we inhabit, we manifestly observe that there are three spatial coordinates, so our space must have at least one dimension more than a two-sphere. It is easy to generalize from the curved two-dimensional manifold (surface) [Equation (2.18)] embedded in three-space to the curved three-dimensional manifold (hypersurface)

$$x^2 + y^2 + z^2 + w^2 = R^2 \quad (2.28)$$

of a three-sphere (hypersphere) embedded in Euclidean four-space with coordinates x , y , z and a fourth fictitious space coordinate w .

Just as the metric in Equation (2.19) could be written without explicit use of z , the metric on the three-sphere in Equation (2.28) can be written without use of w ,

$$dl^2 = dx^2 + dy^2 + dz^2 + \frac{(x dx + y dy + z dz)^2}{R^2 - x^2 - y^2 - z^2}, \quad (2.29)$$

or, in the more convenient spherical coordinates used in Equation (2.27),

$$dl^2 = R^2(t) \left(\frac{R^2 d\sigma^2}{R^2 - (R\sigma)^2} + \sigma^2 d\theta^2 + \sigma^2 \sin^2\theta d\phi^2 \right). \quad (2.30)$$

Note that the introduction of the comoving coordinate σ in Equation (2.27) does not affect the parameter R defining the hypersurface in Equation (2.28). No point is preferred on the manifold [Equation (2.28)], and hence it can describe a spatially homogeneous and isotropic three-dimensional universe in accord with the cosmological principle.

Another example of a curved Riemannian two-space is the surface of a hyperboloid obtained by changing the sign of R^2 in Equation (2.18). The geodesics are hyperbolas, the surface is also unbounded, but in contrast to the spherical surface it is infinite in extent. It can also be generalized to a three-dimensional curved hypersurface, a three-hyperboloid, defined by Equation (2.28) with R^2 replaced by $-R^2$.

The Gaussian curvature of all geodesic three-surfaces in Euclidean four-space is

$$K = k/R^2, \quad (2.31)$$

where the *curvature parameter* k can take the values $+1$, 0 , -1 , corresponding to the three-sphere, flat three-space, and the three-hyperboloid, respectively. Actually, k can take any positive or negative value, because we can always rescale σ to take account of different values of k .

The Robertson–Walker Metric. Let us now include the time coordinate t , the cosmic expansion scale $a(t)$ and the curvature parameter k in Equation (2.29). We then obtain the complete metric derived by *Howard Robertson* and *Arthur Walker* in 1934:

$$ds^2 = c^2 dt^2 - dl^2 = c^2 dt^2 - a(t)^2 \left(\frac{d\sigma^2}{1 - k\sigma^2} + \sigma^2 d\theta^2 + \sigma^2 \sin^2\theta d\phi^2 \right). \quad (2.32)$$

In the tensor notation of Equation (2.14) the components of the Robertson–Walker metric \mathbf{g} are obviously

$$g_{00} = 1, \quad g_{11} = -\frac{a^2(t)}{1 - k\sigma^2}, \quad g_{22} = -a^2(t)\sigma^2, \quad g_{33} = -a^2(t)\sigma^2 \sin^2\theta. \quad (2.33)$$

Equation (2.32) is the most general metric for a four-dimensional space-time which is homogeneous and isotropic at a given time. It then *will always remain homogeneous and isotropic*, because a galaxy at the point (σ, θ, ϕ) will always remain at that point, only the scale of spatial distances $a(t)$ changing with time. The displacements will be $d\sigma = d\theta = d\phi = 0$ and the metric equation will reduce to

$$ds^2 = c^2 dt^2. \quad (2.34)$$

For this reason one calls such an expanding frame a *comoving frame*. A metric with $k = 0$ is called *flat*.

An observer at rest in the comoving frame is called a *fundamental observer*. If the Universe appears to be homogeneous to him/her, it must also be isotropic. But another observer located at the same point and in relative motion with respect to the fundamental observer does not see the Universe as isotropic. Thus the comoving frame is really a preferred frame, and a very convenient one, as we shall see later in conjunction with the cosmic background radiation. Let us note here that a fundamental observer may find that not all astronomical bodies recede radially; a body at motion relative to the comoving coordinates (σ, θ, ϕ) will exhibit peculiar motion in other directions.

Another convenient comoving coordinate is χ , defined by integrating over

$$d\chi = \frac{d\sigma}{\sqrt{1 - k\sigma^2}}. \quad (2.35)$$

Inserting this into Equation (2.32), the metric can be written

$$ds^2 = c^2 dt^2 - a^2(t)[d\chi^2 + S_k^2(\chi)(d\theta^2 + \sin^2\theta d\phi^2)], \quad (2.36)$$

where

$$S_k(\chi) \equiv \sigma$$

and

$$S_1(\chi) = \sin \chi, \quad S_0(\chi) = \chi, \quad S_{-1}(\chi) = \sinh \chi. \quad (2.37)$$

We shall use the metrics in Equations (2.32) and (2.36) interchangeably since both offer advantages. In so doing I have take great care of not introducing contradictions.

Let us briefly digress to define what is sometimes called *cosmic time*. In an expanding universe the galaxies are all moving away from each other (let us ignore peculiar

velocities) with clocks running at different local time, but from our vantage point we would like to have a time value applicable to all of them. If one postulates, with *Hermann Weyl* (1885–1955), that the expansion is so regular that the world lines of the galaxies form a nonintersecting and diverging three-bundle of geodesics, and that one can define spacelike hypersurfaces which are orthogonal to all of them. Then each such hypersurface can be labeled by a constant value of the time coordinate x^0 , and using this value one can meaningfully talk about cosmic epochs for the totality of the Universe. This construction in space-time does not imply the choice of a preferred time in conflict with special relativity.

2.3 Relativistic Distance Measures

Let us consider how to measure distances in our comoving frame in which we are at the origin. The *comoving distance* from us to a galaxy at comoving coordinates $(\sigma, 0, 0)$ is not an observable because a distant galaxy can only be observed by the light it emitted at an earlier time $t < t_0$. In a space-time described by the Robertson–Walker metric the light signal propagates along the geodesic $ds^2 = 0$. Choosing $d\theta^2 = d\phi^2 = 0$, it follows from Equation (2.36) that this geodesic is defined by

$$c^2 dt^2 - a^2(t) d\chi^2 = 0.$$

It follows that χ can be written

$$\chi = c \int_t^{t_0} \frac{dt}{a(t)}. \quad (2.38)$$

The time integral in Equation (2.38) is called the *conformal time*.

Proper Distance. Let us now define the *proper distance* d_p at time t_0 (when the cosmic scale is $a(t_0) = 1$) to the galaxy at $(\sigma, 0, 0)$. This is a function of σ and of the intrinsic geometry of space-time and the value of k . Integrating the spatial distance $dl \equiv |d\mathbf{l}|$ in Equation (2.32) from 0 to d_p we find

$$d_p = \int_0^\sigma \frac{d\sigma}{\sqrt{1 - k\sigma^2}} = \frac{1}{\sqrt{k}} \sin^{-1}(\sqrt{k}\sigma) = \chi. \quad (2.39)$$

For flat space $k = 0$ we find the expected result $d_p = \sigma$. In a universe with curvature $k = +1$ and scale a then Equation (2.39) becomes

$$d_p = a\chi = a \sin^{-1}\sigma \quad \text{or} \quad \sigma = \sin(d_p/a).$$

As the distance d_p increases from 0 to $\frac{1}{2}\pi a$, σ also increases from 0 to its maximum value 1. However, when d_p increases from $\frac{1}{2}\pi a$ to πa , σ decreases back to 0. Thus, travelling a distance $d_p = \pi a$ through the curved three-space brings us to the other end of the Universe. Travelling on from $d_p = \pi a$ to $d_p = 2\pi a$ brings us back to the point of departure. In this sense a universe with positive curvature is closed.

Similarly, the area of a three-sphere centered at the origin and going through the galaxy at σ is

$$A = 4\pi a^2 \sigma^2 = 4\pi a^2 \sin^2(d_p/a). \quad (2.40)$$

Clearly, A goes through a maximum when $d_p = \frac{1}{2}\pi a$, and decreases back to 0 when d_p reaches πa . Note that $A/4$ equals the area enclosed by the circle formed by intersecting a two-sphere of radius R with a horizontal plane, as shown in Figure 2.1. The intersection with an equatorial plane results in the circle enclosing maximal area, $A/4 = \pi R^2$, all other intersections making smaller circles. A plane tangential at either pole has no intersection, thus the corresponding ‘circle’ has zero area.

The volume of the three-sphere in Equation (2.28) can then be written in analogy with Equation (2.23),

$$V = 2 \int_0^{2\pi} d\phi \int_0^\pi d\theta \int_0^1 d\sigma \sqrt{\det \mathbf{g}_{\text{RW}}}, \quad (2.41)$$

where the determinant of the spatial part of the Robertson–Walker metric matrix \mathbf{g}_{RW} is now

$$\det \mathbf{g}_{\text{RW}} = a^6 \frac{\sigma^4}{1 - \sigma^2} \sin^2 \theta. \quad (2.42)$$

The factor 2 in Equation (2.41) comes from the sign ambiguity of w in Equation (2.28). Both signs represent a complete solution. Inserting Equation (2.42) into Equation (2.41) one finds the volume of the three-sphere:

$$V = 2\pi^2 a^3. \quad (2.43)$$

The hyperbolic case is different. Setting $k = -1$, the function in Equation (2.39) is $i^{-1} \sin^{-1} i\sigma \equiv \sinh^{-1} \sigma$, thus

$$d_p = a\chi = a \sinh^{-1} \sigma \quad \text{or} \quad \sigma = \sinh(d_p/a). \quad (2.44)$$

Clearly this space is open because σ grows indefinitely with d_p . The area of the three-hyperboloid through the galaxy at σ is

$$A = 4\pi a^2 \sigma^2 = 4\pi a^2 \sinh^2(d_p/a). \quad (2.45)$$

Let us differentiate d_p in Equation (2.39) with respect to time, noting that σ is a constant since it is a comoving coordinate. We then obtain the Hubble flow v experienced by a galaxy at distance d_p :

$$v = \dot{d}_p = \dot{a}(t) \int_0^\sigma \frac{d\sigma}{\sqrt{1 - k\sigma^2}} = \frac{\dot{a}(t)}{a(t)} d_p. \quad (2.46)$$

Thus the Hubble flow is proportional to distance, and Hubble’s law emerges in a form more general than Equation (1.20):

$$H(t) = \frac{\dot{a}(t)}{a(t)}. \quad (2.47)$$

Recall that v is the velocity of expansion of the space-time geometry. A galaxy with zero comoving velocity would appear to have a radial recession velocity v because of the expansion.

Particle and Event Horizons. In Equation (1.14) we defined the Hubble radius r_H as the distance reached in one Hubble time, τ_H , by a light signal propagating along a straight line in flat, static space. Let us define the *particle horizon* σ_{ph} or χ_{ph} (also *object horizon*) as the largest comoving spatial distance from which a light signal could have reached us if it was emitted at time $t = t_{\text{min}} < t_0$. Thus it delimits the size of that part of the Universe that has come into causal contact since time t_{min} . If the past time t is set equal to the last scattering time (the time when the Universe became transparent to light, and thus the earliest time anything was visible, as we will discuss in a later chapter) the particle horizon delimits the visible Universe. From Equation (2.38),

$$\chi_{\text{ph}} = c \int_{t_{\text{min}}}^{t_0} \frac{dt}{a(t)}, \quad (2.48)$$

and from the notation in Equation (2.37),

$$\sigma_{\text{ph}} = S_k(\chi_{\text{ph}}).$$

A particle horizon exists if t_{min} is in the finite past. Clearly the value of σ_{ph} depends sensitively on the behavior of the scale of the Universe at that time, $a(t_{\text{ph}})$.

If $k \geq 0$, the proper distance (subscript ‘P’) to the particle horizon (subscript ‘ph’) at time t is

$$d_{\text{P,ph}} = a(t)\chi_{\text{ph}}. \quad (2.49)$$

Note that d_{P} equals the Hubble radius $r_H = c/H_0$ when $k = 0$ and the scale is a constant, $a(t) = a$. When $k = -1$ the Universe is open, and $d_{\text{P,ph}}$ cannot be interpreted as a measure of its size.

In an analogous way, the comoving distance σ_{eh} to the *event horizon* is defined as the spatially most distant present event from which a world line can ever reach our world line. By ‘ever’ we mean a finite future time, t_{max} :

$$\chi_{\text{eh}} \equiv c \int_{t_0}^{t_{\text{max}}} \frac{dt}{a(t)}. \quad (2.50)$$

The particle horizon σ_{ph} at time t_0 lies on our past light cone, but with time our particle horizon will broaden so that the light cone at t_0 will move inside the light cone at $t > t_0$ (see Figure 2.1). The event horizon at this moment can only be specified given the time distance to the ultimate future, t_{max} . Only at t_{max} will our past light cone encompass the present event horizon. Thus the event horizon is our ultimate particle horizon. Comoving bodies at the Hubble radius recede with velocity c , but the particle horizon itself recedes even faster. From

$$d(Hd_{\text{P,ph}})/dt = \dot{H}d_{\text{P,ph}} + Hd_{\text{P,ph}} = 0,$$

and making use of the *deceleration parameter* q , defined by

$$q = -\frac{a\ddot{a}}{\dot{a}^2} = -\frac{\ddot{a}}{aH^2}, \quad (2.51)$$

one finds

$$\dot{d}_{\text{P,ph}} = c(q + 1). \quad (2.52)$$

Thus when the particle horizon grows with time, bodies which were at spacelike distances at earlier times enter into the light cone.

The integrands in Equations (2.48) and (2.50) are obviously the same, only the integration limits are different, showing that the two horizons correspond to different conformal times. If $t_{\min} = 0$, the integral in Equation (2.48) may well diverge, in which case there is no particle horizon. Depending on the future behavior of $a(t)$, an event horizon may or may not exist. If the integral diverges as $t \rightarrow \infty$, every event will sooner or later enter the event horizon. The event horizon is then a function of waiting time only, but there exists no event horizon at $t = \infty$. But if $a(t)$ accelerates, so that distant parts of the Universe recede faster than light, then there will be an ultimate event horizon. We shall see later that $a(t)$ indeed appears to accelerate.

Redshift and Proper Distance. In Equation (1.20) in the previous chapter we parametrized the rate of expansion \dot{a} by the Hubble constant H_0 . It actually appeared as a dynamical parameter in the lowest-order Taylor expansion of $a(t)$, Equation (1.18). If we allow $H(t)$ to have some mild time dependence, that would correspond to introducing another dynamical parameter along with the next term in the Taylor expansion. Thus adding the second-order term to Equation (1.18), we have for $a(t)$,

$$a(t) \approx 1 - \dot{a}_0(t_0 - t) + \frac{1}{2}\ddot{a}_0(t_0 - t)^2. \quad (2.53)$$

Making use of the definition in Equation (2.47), the second-order expansion for the dimensionless scale factor is

$$a(t) \approx 1 - H_0(t_0 - t) + \frac{1}{2}\dot{H}_0(t_0 - t)^2. \quad (2.54)$$

As long as the observational information is limited to the first time derivative \dot{a}_0 , no further terms can be added to these expansions. To account for \ddot{a}_0 , we shall now make use of the present value of the deceleration parameter in Equation (2.51), q_0 . Then the lowest-order expression for the cosmological redshift, Equation (1.18), can be replaced by

$$\begin{aligned} z(t) &= (a(t))^{-1} - 1 \\ &= \left[1 - H_0(t - t_0) - \frac{1}{2}q_0 H_0^2(t - t_0)^2 \right]^{-1} - 1 \\ &\approx -H_0(t - t_0) + \left(1 + \frac{1}{2}q_0 \right) H_0^2(t - t_0)^2. \end{aligned}$$

This expression can further be inverted to express $H_0(t - t_0)$ as a function of the redshift to second order:

$$H_0(t_0 - t) \approx z - \left(1 + \frac{1}{2}q_0 \right) z^2. \quad (2.55)$$

Let us now find the proper distance d_p to an object at redshift z in this approximation. Eliminating χ in Equations (2.38) and (2.39) we have

$$d_p = c \int_t^{t_0} \frac{dt}{a(t)}.$$

We then insert $a(t)$ from Equation (2.54) to lowest order in $t_0 - t$, obtaining

$$d_P \approx c \int_t^{t_0} [1 + H_0(t_0 - t)] dt = c(t_0 - t) \left[1 + \frac{1}{2} H_0(t_0 - t) \right]. \quad (2.56)$$

Substituting the expression in Equation (2.55) into this yields the sought result:

$$d_P(z) \approx \frac{c}{H_0} \left(z - \frac{1}{2}(1 + q_0)z^2 \right). \quad (2.57)$$

The first term on the right gives Hubble's linear law [Equation (1.15)], and thus the second term measures deviations from linearity to lowest order. The parameter value $q_0 = -1$ obviously corresponds to no deviation. The linear law has been used to determine H_0 from galaxies within the Local Supercluster (LSC). On the other hand, one also observes deceleration of the expansion in the local universe due to the lumpiness of matter. For instance, the local group clearly feels the overdensity of the Virgo cluster at a distance of about 17 Mpc, falling towards it with a peculiar velocity of about 630 km s^{-1} [3]. It has been argued that the peculiar velocities in the LSC cannot be understood without the pull of the neighboring Hydra–Centaurus supercluster and perhaps a still larger overdensity in the supergalactic plane, a rich cluster (the A3627) nicknamed the ‘Great Attractor’.

It should be clear from this that one needs to go to even greater distances, beyond the influences of local overdensities, to determine a value for q_0 . Within the LSC it is safe to conclude that only the linear term in Hubble's law is necessary.

Equation (2.57) is the conventional formula, which is a good approximation for small z . The approximation obviously deteriorates as z increases, so that it attains its maximum at $z = 1/1 + q_0$. In Figure 2.5 we plot the function $d_P(z)$ for small values of z .

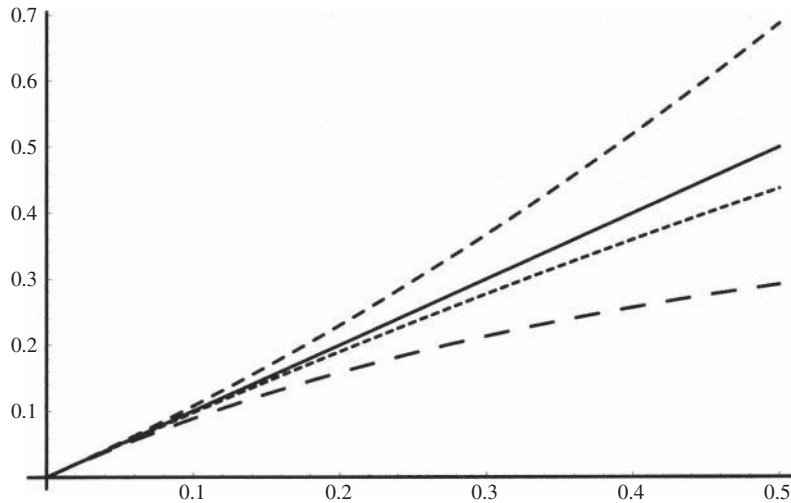


Figure 2.5 Approximate distance measures d to second order in z . The solid curve shows the linear function, Equation (1.15); the short-dashed curve the proper distance d_P , Equation (2.57); the medium-dashed curve the luminosity distance d_L , Equation (2.60); and the long-dashed curve the angular size distance d_A , Equation (2.64). The value of q_0 is -0.5 .

Redshift and Luminosity Distance. Consider an astronomical object emitting photons isotropically with power or absolute luminosity L . At the *luminosity distance* d_L from the object we observe only the fraction B_s , its surface brightness, given by the inverse-square distance law

$$B_s = \frac{L}{4\pi d_L^2}. \quad (2.58)$$

Let us now find d_L as a function of z in such a way that the Euclidean inverse-square law [Equation (2.58)] is preserved. If the Universe does not expand and the object is stationary at proper distance d_p , a telescope with area A will receive a fraction $A/4\pi d_p^2$ of the photons. But in a universe characterized by an expansion $a(t)$, the object is not stationary, so the energy of photons emitted at time t_e is redshifted by the factor $(1+z) = a^{-1}(t_e)$. Moreover, the arrival rate of the photons suffers time dilation by another factor $(1+z)$, often called the *energy effect*. The *apparent brightness* B_a is then given by

$$B_a = \frac{L}{4\pi d_p^2(1+z)^2}. \quad (2.59)$$

Equating $B_a = B_s$ one sees that $d_L = d_p(1+z)$, and making use of the expression in Equation (2.57) one obtains

$$d_L(z) \approx \frac{c}{H_0} \left(z + \frac{1}{2}(1-q_0)z^2 \right). \quad (2.60)$$

In Figure 2.5 we plot the function $d_L(z)$ for small values of z .

Astronomers usually replace L and B by two empirically defined quantities, *absolute magnitude* M of a luminous object and *apparent magnitude* m . The replacement rule is

$$m - M = -5 + 5 \log d_L, \quad (2.61)$$

where d_L is expressed in parsecs (pc) and the logarithm is to base 10. For example, if one knows the distance d_L to a galaxy hosting a supernova, its absolute magnitude M can be obtained from observations of its apparent magnitude m .

Parallax Distance. Some measurements of H_0 depend directly on the calibration of local distance indicators which form the first rung of a ladder of distance measurements. The distances to relatively nearby stars can be measured by the *trigonometrical parallax* up to about 30 pc away (see Table A.1 in the appendix for cosmic distances). This is the difference in angular position of a star as seen from Earth when at opposite points in its circumsolar orbit. The *parallax distance* d_p is defined as

$$d_p = d_p / \sqrt{1 - k\sigma^2}. \quad (2.62)$$

It has been found that most stars in the Galaxy for which we know the luminosity from a kinematic distance determination exhibit a relationship between surface temperature or color and absolute luminosity, the *Hertzsprung–Russell* relation. These stars are called *main-sequence stars* and they sit on a fairly well-defined curve in the

temperature–luminosity plot. Temperature can be determined from color—note that astronomers define color as the logarithm of the ratio of the apparent brightnesses in the red and the blue wavelength bands. Cool stars with surface temperature around 3000 K are infrared, thus the part of their spectrum which is in the visible is dominantly red. Hot stars with surface temperature around 12 000 K are ultraviolet, thus the part of their spectrum which is in the visible is dominantly blue. The Sun, with a surface temperature of 5700 K, radiates mainly in the visible, thus its color is a blended white, slightly yellow. Most main-sequence stars like our Sun are in a prolonged state of steady burning of hydrogen into helium.

Once this empirical temperature–luminosity relation is established, it can be used the other way around to derive distances to farther main-sequence stars: from their color one obtains the luminosity which subsequently determines d_L . By this method one gets a second rung in a ladder of estimates which covers distances within our Galaxy.

Angular Size Distance. Yet another measure of distance is the *angular size distance* d_A . In Euclidean space an object of size D that is at distance d_A will subtend an angle 2θ such that

$$2\theta = \tan(D/d_A) \approx D/d_A,$$

where the approximation is good for small θ . This can serve as the definition of d_A in Euclidean space. In general relativity we can still use this equation to define a distance measure d_A . From the metric in Equation (2.32) the radius of a source of light at comoving distance σ is $D = a\sigma\theta$, so

$$d_A = D/\theta = a\sigma = aS_k d_P. \quad (2.63)$$

This definition preserves the relation between angular size and distance, a property of Euclidean space. But expansion of the Universe and the changing scale factor $a(t)$ means that as proper distance d_P or redshift z increases, the angular diameter distance initially increases but ultimately decreases. Light rays from the object detected by the observer have been emitted when the proper distance to the object, measured at fixed world time, was small. Because the proper distance between observer and source is increasing faster than the speed of light, emitted light in the direction of the observer is initially moving away from the observer.

The redshift dependence of d_A can be found from Equations (2.57) and (2.37) once k is known. In Figure 2.5 we plot d_A for the choice $k = 0$ when

$$d_A = ad_P = \frac{d_P}{1+z}. \quad (2.64)$$

The k dependence makes it a useful quantity to determine cosmological parameters. In particular, k is sensitive to certain combinations of well-measured parameters.

Distance Ladder Continued. As the next step on the distance ladder one chooses calibrators which are stars or astronomical systems with specific uniform properties, so called *standard candles*. The *RR Lyrae* stars all have similar absolute luminosities,

and they are bright enough to be seen out to about 300 kpc. A very important class of standard candles are the *Cepheid* stars, whose absolute luminosity oscillates with a constant period P in such a way that $\log P \propto 1.3 \log L$. The period P can be observed with good precision, thus one obtains a value for L . Cepheids have been found within our Galaxy where the period–luminosity relation can be calibrated by distances from trigonometric parallax measurements. This permits use of the period–luminosity relation for distances to Cepheids within the *Large Magellanic Cloud* (LMC), our satellite galaxy. At a distance of 55 kpc the LMC is the first important extragalactic landmark.

The resolution of individual stars within galaxies clearly depends on the distance to the galaxy. This method, called *surface-brightness fluctuations* (SBFs), is an indicator of relative distances to elliptical galaxies and some spirals. The internal precision of the method is very high, but it can be applied only out to about 70 Mpc.

Globular clusters are gravitationally bound systems of 10^5 – 10^6 stars forming a spherical population orbiting the center of our Galaxy. From their composition one concludes that they were created very early in the evolution of the Galaxy. We already made use of their ages to estimate the age of the Universe in Section 1.5. Globular clusters can also be seen in many other galaxies, and they are visible out to 100 Mpc. Within the Galaxy their distance can be measured as described above, and one then turns to study the statistical properties of the clusters: the frequency of stars of a given luminosity, the mean luminosity, the maximum luminosity, and so on. A well-measured cluster then becomes a standard candle with properties presumably shared by similar clusters at all distances. Similar statistical indicators can be used to calibrate clusters of galaxies; in particular the brightest galaxy in a cluster is a standard candle useful out to 1 Gpc.

The next two important landmarks are the distances to the Virgo cluster, which is the closest moderately rich concentrations of galaxies, and to the Coma cluster, which is one of the closest clusters of high richness. The Virgo distance has been determined to be 17 Mpc by the observations of galaxies containing several Cepheids, by the *Hubble Space Telescope* [4]. The Coma is far enough, about 100 Mpc, that its redshift is almost entirely due to the cosmological expansion.

The existence of different methods of calibration covering similar distances is a great help in achieving higher precision. The expansion can be verified by measuring the surface brightness of standard candles at varying redshifts, the *Tolman test*. If the Universe does indeed expand, the intensity of the photon signal at the detector is further reduced by a factor $(1+z)^2$ due to an optical aberration which makes the surface area of the source appear increased. Such tests have been done and they do confirm the expansion.

The *Tully–Fisher* relation is a very important tool at distances which overlap those calibrated by Cepheids, globular clusters, galaxy clusters and several other methods. This empirical relation expresses correlations between intrinsic properties of whole spiral galaxies. It is observed that their absolute luminosity and their circular rotation velocity v_c are related by

$$L \propto v_c^4. \quad (2.65)$$

The Tully–Fisher relation for spiral galaxies is calibrated by nearby spiral galaxies having Cepheid calibrations, and it can then be applied to spiral galaxies out to 150 Mpc. Elliptical galaxies do not rotate, they are found to occupy a *fundamental plane* in which an effective radius is tightly correlated with the surface brightness inside that radius and with the central velocity dispersion of the stars. In principle, this method could be applied out to $z \approx 1$, but in practice stellar evolution effects and the nonlinearity of Hubble’s law limit the method to $z \lesssim 0.1$, or about 400 Mpc.

For more details on distance measurements the reader is referred to the excellent treatment in the book by Peacock [5].

2.4 Tests of Special Relativity

Special relativity is an inseparable part of quantum field theory which describes the world of elementary particles with an almost incredible precision. Particle physics has tested special relativity in thousands of different experiments without finding a flaw: the Lorentz invariance is locally exact. But at astronomical and cosmological scales the local Lorentz invariance has to be replaced by General Relativity. The quibble about whether special relativity is generally true and testable at cosmological distances and time scales is therefore meaningless.

Special relativity really contains only one parameter, c , the velocity of light *in vacuo*, which has the dimension of length/time. Could c be variable, or is c even measurable at all? One is free to choose $c = 1$ locally because that only implies a rescaling of the units of length. To explain problematic observations the possibility of a *variable speed of light* has sometimes been invoked. Also Newton’s constant of gravitation, G , could in principle be variable. Already Einstein admitted that the value of G could depend on the local strength of the gravitational field. In order to avoid trivial rescaling of units, one must test the simultaneous variation of c , G and the fine structure constant which can be combined to a dimensionless number. Note that c and G enter in the combination G/c^4 in the Einstein Equation (3.29).

A measurement of the radius of Mercury has produced no time variation of c

$$\frac{\dot{c}}{c} = 0 \pm 2 \times 10^{-12} \text{ yr}^{-1}. \quad (2.66)$$

The white dwarf star Stein 2051B within our Galaxy also provides only a limit when combined with the upper limit of the variability of the fine structure constant,

$$H_0 \frac{\dot{c}}{c} = 0 \pm 2.1 \times 10^{-3} \text{ yr} \quad (2.67)$$

using $H_0 = 1.5 \times 10^{10} \text{ yr}$.

Determinations of the time variation of G have also given null results. The strongest constraint due to a lunar laser ranging experiment gives

$$\left| \frac{\dot{G}}{G} \right| \leq 1.3 \times 10^{-12} \text{ yr}^{-1}. \quad (2.68)$$

Variable speed of light is in conflict with Lorentz invariance and with Einstein's time dilation formula in Equation (2.5). Although there have been many attempts to test this by measuring the transverse second-order Doppler shift by Mössbauer spectroscopy, the results are claimed to be either wrong or doubtful, and in need of improved technology.

On the theoretical side there are interesting generalizations of the linear Lorentz transformations to uniformly accelerated or rotating frames. Some generalized transformations predict acceleration-dependent Doppler shift and time dilation, as well as a maximal acceleration [6, 7].

Problems

1. Starting from the postulates (i) and (ii) in Section 2.1 and the requirement that ds^2 in Equation (2.1) should have the same form in the primed and unprimed coordinates, derive the linear Lorentz transformation in Equation (2.2) and the expression in Equation (2.3).
2. The radius of the Galaxy is 3×10^{20} m. How fast would a spaceship have to travel to cross it in 300 yr as measured on board? Express your result in terms of $\gamma = 1/\sqrt{1 - v^2/c^2}$ [7].
3. An observer sees a spaceship coming from the west at a speed of $0.6c$ and a spaceship coming from the east at a speed $0.8c$. The western spaceship sends a signal with a frequency of 10^4 Hz in its rest frame. What is the frequency of the signal as perceived by the observer? If the observer sends on the signal immediately upon reception, what is the frequency with which the eastern spaceship receives the signal [7]?
4. If the eastern spaceship in the previous problem were to interpret the signal as one that is Doppler shifted because of the relative velocity between the western and eastern spaceships, what would the eastern spaceship conclude about the relative velocity? Show that the relative velocity must be $(v_1 + v_2)/(1 + v_1 v_2/c^2)$, where v_1 and v_2 are the velocities as seen by an outside observer [7].
5. A source flashes with a frequency of 10^{15} Hz. The signal is reflected by a mirror moving away from the source with speed 10 km s^{-1} . What is the frequency of the reflected radiation as observed at the source [7]?
6. Suppose that the evolution of the Universe is described by a constant decelerating parameter $q = \frac{1}{2}$. We observe two galaxies located in opposite directions, both at proper distance d_p . What is the maximum separation between the galaxies at which they are still causally connected? Express your result as a fraction of distance to d_p . What is the observer's particle horizon?
7. Show that the Hubble distance $r_H = c/H$ recedes with radial velocity

$$\dot{r}_H = c(1 + q). \quad (2.69)$$

8. Is the sphere defined by the Hubble radius r_H inside or outside the particle horizon?

9. Calculate whether the following space-time intervals from the origin are space-like, timelike or lightlike: $(1, 3, 0, 0)$; $(3, 3, 0, 0)$; $(3, -3, 0, 0)$; $(0, 3, 0, 0)$; $(3, 1, 0, 0)$ [1].
10. The supernova 1987A explosion in the Large Magellanic Cloud 170 000 light years from Earth produced a burst of anti-neutrinos $\bar{\nu}_e$ which were observed in terrestrial detectors. If the anti-neutrinos are massive, their velocity would depend on their mass as well as their energy. What is the proper time interval between the emission, assumed to be instantaneous, and the arrival on Earth? Show that in the limit of vanishing mass the proper time interval is zero. What information can be derived about the anti-neutrino mass from the observation that the energies of the anti-neutrinos ranged from 7 to 11 MeV, and the arrival times showed a dispersion of 7 s?
11. The theoretical resolving power of a telescope is given by $\alpha = 1.22\lambda/D$, where λ is the wavelength of the incoming light and D is the diameter of the mirror. Assuming $D = 5$ m and $\lambda = 8 \times 10^{-7}$ m, determine the largest distance to a star that can be measured by the parallax method. (In reality, atmospheric disturbances set tighter limits.)

References

- [1] Kenyon, I. R. 1990 *General relativity*. Oxford University Press, Oxford.
- [2] Pyykkö, P. 1988 *Chem. Rev.* **88**, 563.
- [3] Lynden-Bell, D. *et al.* 1988 *Astrophys. J.* **326**, 19.
- [4] Freedman, W. L. *et al.* 1994 *Nature* **371**, 757.
- [5] Peacock, J. A. 1999 *Cosmological physics*. Cambridge University Press, Cambridge.
- [6] Friedman, Y. and Scarr, T. 2013 *Physica Scripta* **87**, 055004, and references therein.
- [7] Gasiörowicz, S. 1979 *The structure of matter*. Addison-Wesley, Reading, MA.

General Relativity

Newton's law of gravitation, Equation (1.28), runs into serious conflict with special relativity in three different ways. First, there is no obvious way of rewriting it in terms of invariants, since it only contains scalars. Second, it has no explicit time dependence, so gravitational effects propagate instantaneously to every location in the Universe, in fact, also infinitely far outside the horizon of the Universe!

Third, the *gravitating mass* m_G appearing in Equation (1.28) is totally independent of the *inert mass* m appearing in Newton's second law [Equation (2.29)], as we already noted, yet for unknown reasons both masses appear to be equal to a precision of 10^{-13} or better (10^{-18} is expected soon). Clearly a theory is needed to establish a formal link between them. Mach thought that the inert mass of a body was somehow linked to the gravitational mass of the whole Universe. To be rigorous, he thought that axes of local nonrotating frames, such as axes of gyroscopes, in their time-evolution precisely follow some average of the motion of matter in the Universe. Einstein, who was strongly influenced by the ideas of Mach, called this *Mach's principle*. In his early work on general relativity he considered it to be one of the basic, underlying principles, together with the principles of equivalence and covariance, but in his later publications he no longer referred to it. This may have been a misunderstanding of Einstein because, if one carefully defines the words *precisely follow* and *some average*, Mach's principle is a consequence of cosmology with Einstein Gravity.

Facing the above shortcomings of Newtonian mechanics and the limitations of special relativity Einstein set out on a long and tedious search for a better law of gravitation valid in the inhomogeneous gravitational field near a massive body, yet one that would reduce to Newton's law in some limit. Realizing that the space we live in was not flat, except locally *equivalent* to (a patchwork of flat frames describing) a curved space-time, the law of gravitation has to be a covariant relation between mass density and curvature. Thus Einstein proceeded to combine the principle of equivalence which we describe in Section 3.1 and the principle of general covariance which we meet in Section 3.2.

This is most conveniently done using tensor notation, briefly presented in Section 3.2, which has the advantage of permitting laws of nature to be written in the same form in all invariant frames.

In Section 3.3 we derive Einstein's law of gravitation starting from an *action principle*, the Einstein-Hilbert action. In Section 3.4 we use simple qualitative arguments to rederive Einstein's law in the weak field limit which is equivalent to Newton's law of gravitation.

3.1 The Principle of Equivalence

Consider the lift in Figure 3.1 moving vertically in a tall tower (it is easy to imagine an lift to be at rest with respect to an outside observer fixed to the tower, whereas the more 'modern' example of a spacecraft is not at rest when we observe it to be geostationary). A passenger in the lift testing the law of gravitation would find that objects dropped to the floor acquire the usual gravitational acceleration g when the lift stands still, or moves with constant speed. However, when the outside observer notes

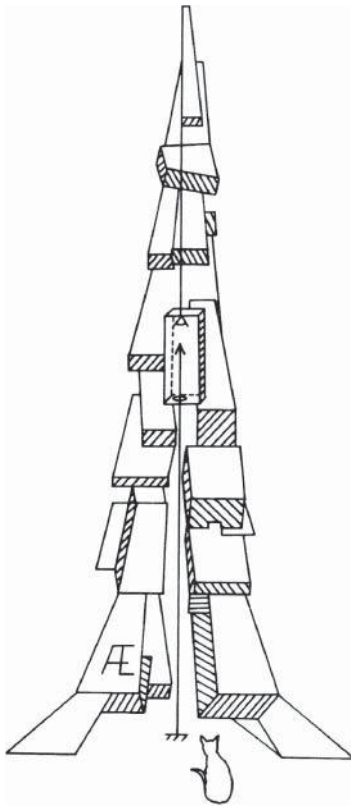


Figure 3.1 The Einstein lift mounted in a nonEuclidean tower. An observer is seen in the foreground.

that the lift is accelerating upwards, tests inside the lift reveal that the objects acquire an acceleration larger than g , and vice versa when the lift is accelerating downwards. In the limit of free fall (unpleasant to the passenger) the objects appear weightless, corresponding to zero acceleration.

Let us now replace the lift with a spacecraft with the engines turned off, located at some neutral point in space where all gravitational pulls cancel or are negligible: a good place is the *Lagrange point*, where the terrestrial and solar gravitational fields cancel. All objects, including the pilot, would appear weightless there.

Now turning on the engines by remote radio control, the spacecraft could be accelerated upwards so that objects on board would acquire an acceleration g towards the floor. The pilot would then rightly conclude that

gravitational pull and local acceleration are equivalent

and indistinguishable if no outside information is available and if $m = m_G$. This conclusion forms the *weak equivalence principle (WEP)*, which states that an observer in a gravitational field will not experience free fall as a gravitational effect, but as being at rest in a locally accelerated frame.

A passenger in the lift measuring g could well decide from his local observations that Earth's gravitation actually does not exist, but that the lift is accelerating radially outwards from Earth. This interpretation does not come into conflict with that of another observer on the opposite side of Earth whose frame would accelerate in the opposite direction, because that frame is only local to him/her.

The weak equivalence principle is already embodied in the *Galilean equivalence principle* in mechanics between motion in a uniform gravitational field and a uniformly accelerated frame of reference. What Einstein did was to generalize this to all of physics, in particular phenomena involving light.

The more general formulation is the important *strong equivalence principle (SEP)*, which states that

to an observer in free fall in a gravitational field the results of all local experiments are completely independent of the magnitude of the field.

In a suitably small lift or spacecraft, curved space-time can always be approximated by flat Minkowski space-time. In the gravitational field of Earth the gravitational acceleration is directed toward its center. Thus the two test bodies in Figure 3.2 with a space-like separation do not actually fall along parallels, but along different radii, so that their separation decreases with time. This phenomenon is called the *tidal effect*, or sometimes the tidal force, since the test bodies move as if an attractive exchange force acted upon them. The classic example is the tide caused by the Moon on the oceans. The force experienced by a body of mass m and diameter d in gravitational interaction with a body of mass M at a distance r is proportional to the differential of the force of attraction [Equation (1.28)] with respect to r . Neglecting the geometrical shapes of the bodies, the tidal force is

$$F_{\text{tidal}} \approx GMmd/r^3.$$

Thus parts of m located at smaller distances r feel a stronger force.

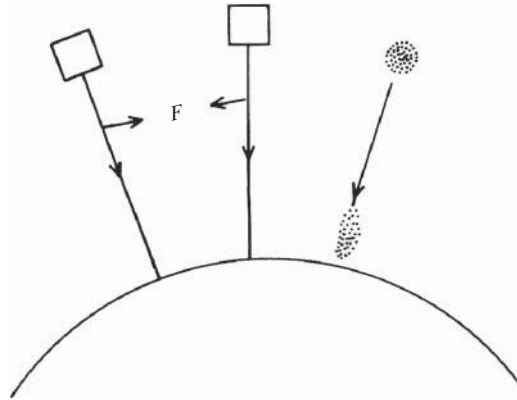


Figure 3.2 Tidal force F acting between two test bodies falling freely towards the surface of a gravitating body. On the right a spherical cluster of small bodies is seen to become ellipsoidal on approaching the body.

An interesting example is offered by a sphere of freely falling particles. Since the strength of the gravitational field increases in the direction of fall, the particles in the front of the sphere will fall faster than those in the rear. At the same time the lateral cross-section of the sphere will shrink due to the tidal effect. As a result, the sphere will be *focused* into an ellipsoid with the same volume. This effect is responsible for the gravitational breakup of very nearby massive stars.

If the tidal effect is too small to be observable, the laboratory can be considered to be local. On a larger scale the gravitational field is clearly quite nonuniform, so if we make use of the equivalence principle to replace this field everywhere by locally flat frames, we get a patchwork of frames which describe a curved space. Since the inhomogeneity of the field is caused by the inhomogeneous distribution of gravitating matter, Einstein realized that the space we live in had to be curved, and the curvature had to be related to the distribution of matter.

But Einstein had already seen the necessity of introducing a four-dimensional space-time, thus it was not enough to describe space-time in a nonuniform gravitational field by a curved space, time also had to be curved. When moving over a patchwork of local and spatially distinct frames, the local time would also have to be adjusted from frame to frame. In each frame the strong equivalence principle requires that measurements of time would be independent of the strength of the gravitational field.

Falling Photons. Let us return once more to the passenger in the Einstein lift for a demonstration of the relation between gravitation and the curvature of space-time. Let the lift be in free fall; the passenger would consider that no gravitational field is present. Standing by one wall and shining a pocket lamp horizontally across the lift, she sees that light travels in a straight path, a geodesic in flat space-time. This is illustrated in Figure 3.3. Thus she concludes that in the absence of a gravitational field space-time is flat.

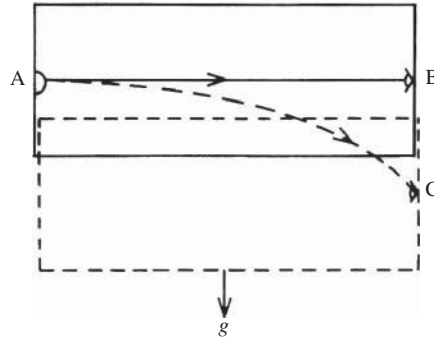


Figure 3.3 A pocket lamp at ‘A’ in the Einstein lift is shining horizontally on a point ‘B’. However, an outside observer who sees that the lift is falling freely with acceleration g concludes that the light ray follows the dashed curve to point ‘C’.

However, the observer outside the tower sees that the lift has accelerated while the light front travels across the lift, and so with respect to the fixed frame of the tower he notices that the light front follows a curved path, as shown in Figure 3.3. He also sees that the lift is falling in the gravitational field of Earth, and so he would conclude that light feels gravitation as if it had mass. He could also phrase it differently: light follows a geodesic, and since this light path is curved it must imply that space-time is curved in the presence of a gravitational field.

When the passenger shines monochromatic light of frequency ν vertically up, it reaches the roof height d in time d/c . In the same time the outside observer records that the lift has accelerated from, say, $v = 0$ to gd/c , where g is the gravitational acceleration on Earth, so that the color of the light has experienced a gravitational redshift by the fraction

$$\frac{\Delta\nu}{\nu} \approx \frac{v}{c} = \frac{gd}{c^2} = \frac{GMd}{r^2c^2}. \quad (3.1)$$

Thus the photons have lost energy ΔE by climbing the distance d against Earth’s gravitational field,

$$\Delta E = h\Delta\nu = -\frac{gdh\nu}{c^2}, \quad (3.2)$$

where h is the *Planck constant*. (Recall that *Max Planck* (1858–1947) was the inventor of the quantization of energy; this led to the discovery and development of *quantum mechanics*.)

If the pocket lamp had been shining electrons of mass m , they would have lost kinetic energy

$$\Delta E = -gmd \quad (3.3)$$

climbing up the distance d . Combining Equations (3.2) and (3.3) we see that the photons appear to possess mass:

$$m = \frac{h\nu}{c^2}. \quad (3.4)$$

Equation (3.1) clearly shows that light emerging from a star with mass M is redshifted in proportion to M . Thus part of the redshift observed is due to this gravitational effect. From this we can anticipate the existence of stars with so large a mass that their gravitational field effectively prohibits radiation to leave. These are the *black holes* to which we shall return later.

Superluminal Photons. A cornerstone in special relativity is that no material particle can be accelerated beyond c , no physical effect can be propagated faster than c , and no signal can be transmitted faster than c . It is an experimental fact that no particle has been found travelling at superluminal speed, but a name for such particles has been invented, *tachyons*. Special relativity does not forbid tachyons, but if they exist they cannot be retarded to speeds below c . In this sense the speed of light constitutes a two-way barrier: an upper limit for ordinary matter and a lower limit for tachyons.

On quantum scales this may be violated, since the photon may appear to possess a mass caused by its interaction with virtual electron–positron pairs. In sufficiently strong curvature fields, the trajectory of a photon may then be distorted through the interaction of gravity on this mass and on the photon’s polarization vector, so that the photon no longer follows its usual geodesic path through curved space-time. The consequence is that SPE may be violated at quantum scales, the photon’s lightcone is changed, and it may propagate with superluminal velocity. This effect, called *gravitational birefringence*, can occur because general relativity is not constructed to obey quantum theory. It may still modify our understanding of the origin of the Universe, when the curvature must have been extreme, and perhaps other similar situations like the interior of black holes. For a more detailed discussion of this effect, see Shore [2] and references therein.

3.2 The Principle of Covariance

Tensors. In four-dimensional space-time all spatial three-vectors have to acquire a zeroth component just like the line element four-vector ds in Equations (2.7) and (2.14). A vector \mathbf{A} with components A_μ in a coordinate system x_μ can be expressed in a transformed coordinate system x'_ν as the vector \mathbf{A}' with components

$$A'_\nu = \frac{\partial x'_\nu}{\partial x_\mu} A_\mu, \quad (3.5)$$

where summation over the repeated index μ is implied, just as in Equation (2.14). A vector which transforms in this way is said to be *contravariant*, which is indicated by the *upper index* for the components A^μ .

A vector \mathbf{B} with components B_μ in a coordinate system x^μ , which transforms in such a way that

$$B'_\nu = \frac{\partial x^\mu}{\partial x'^\nu} B_\mu, \quad (3.6)$$

is called *covariant*. This is indicated by writing its components with a *lower index*. Examples of covariant vectors are the tangent vector to a curve, the normal to a surface, and the four-gradient of a four-scalar ϕ , $\partial\phi/\partial x^\mu$.

In general, tensors can have several contravariant and covariant indices running over the dimensions of a manifold. In a d -dimensional manifold a tensor with r indices is of *rank* r and has d^r components. In particular, an $r = 1$ tensor is a vector, and $r = 0$ corresponds to a scalar. An example of a tensor is the assembly of the n^2 components $X^\mu Y^\nu$ formed as the products (not the scalar product!) of the n components of the vector X^μ with the n components of the vector Y^ν . We have already met the rank 2 tensors $\eta_{\mu\nu}$ with components given by Equation (2.14), and the metric tensor $g_{\mu\nu}$.

Any contravariant vector \mathbf{A} with components A^μ can be converted into a covariant vector by the operation

$$A_\nu = g_{\mu\nu} A^\mu.$$

The contravariant metric tensor $g^{\mu\nu}$ is the matrix inverse of the covariant $g_{\mu\nu}$:

$$g_{\sigma\mu} g^{\mu\nu} = \delta_\sigma^\nu. \quad (3.7)$$

The upper and lower indices of any tensor can be lowered and raised, respectively, by operating with $g_{\mu\nu}$ or $g^{\mu\nu}$ and summing over repeated indices. Thus a covariant vector \mathbf{A} with components A_μ can be converted into a contravariant vector by the operation

$$A^\nu = g^{\mu\nu} A_\mu,$$

For a point particle with mass m and total energy

$$E = \gamma mc^2, \quad (3.8)$$

according to Einstein's famous relation, one assigns a momentum four-vector \mathbf{P} with components $p^0 = E/c$, $p^1 = p_x$, $p^2 = p_y$, $p^3 = p_z$, so that E and the linear momentum $\mathbf{p} = m\mathbf{v}$ become two aspects of the same entity, $\mathbf{P} = (E/c, \mathbf{p})$.

The scalar product P^2 is an invariant related to the mass,

$$P^2 = \eta_{\mu\nu} P^\mu P^\nu = \frac{E^2}{c^2} - p^2 = (\gamma mc)^2, \quad (3.9)$$

where $p^2 \equiv |\gamma\mathbf{p}|^2$. For a massless particle like the photon, it follows that the energy equals the three-momentum times c .

Newton's second law in its nonrelativistic form,

$$\mathbf{F} = m\mathbf{a} = m\dot{\mathbf{v}} = \dot{\mathbf{p}}, \quad (3.10)$$

is replaced by the relativistic expression

$$\mathbf{F} = \frac{d\mathbf{P}}{d\tau} = \gamma \frac{d\mathbf{P}}{dt} = \gamma \left(\frac{dE}{c dt}, \frac{d\mathbf{p}}{dt} \right). \quad (3.11)$$

General Covariance. Although Newton's second law Equation (3.11) is invariant under special relativity in any inertial frame, it is not invariant in accelerated frames. Since this law explicitly involves acceleration, special relativity has to be generalized somehow, so that observers in accelerated frames can agree on the value of acceleration. Thus the next necessary step is to search for quantities which remain invariant under an arbitrary acceleration and to formulate the laws of physics in terms of these. Such a formulation is called *generally covariant*. In a curved space-time

described by the Robertson–Walker metric the approach to general covariance is to find appropriate invariants in terms of tensors which have the desired properties.

Since vectors are rank 1 tensors, vector equations may already be covariant. However, dynamical laws contain many other quantities that are not tensors, in particular space-time derivatives such as $d/d\tau$ in Equation (3.11). Space-time derivatives are not invariants because they imply transporting ds along some curve and that makes them coordinate dependent. Therefore we have to start by redefining derivatives and replacing them with new *covariant derivatives*, which are tensor quantities.

To make the space-time derivative of a vector generally covariant one has to take into account that the direction of a parallel-transported vector changes in terms of the local coordinates along the curve as shown in Figure 3.4. The change is certainly some function of the space-time derivatives of the curved space that is described by the metric tensor.

The covariant derivative operator with respect to the proper time τ is denoted $D/D\tau$ (for a detailed derivation see, e.g., references [1] and [3]). Operating with it on the momentum four-vector P^μ results in another four-vector:

$$F^\mu = \frac{DP^\mu}{D\tau} \equiv \frac{dP^\mu}{d\tau} + \Gamma_{\sigma\nu}^\mu P^\sigma \frac{dx^\nu}{d\tau}. \quad (3.12)$$

The second term contains the changes this vector undergoes when it is parallel transported an infinitesimal distance $c d\tau$. The quantities $\Gamma_{\sigma\nu}^\mu$, called *affine connections*, are readily derivable functions of the derivatives of the metric $g_{\mu\nu}$ in curved space-time, but they are not tensors. Their form is

$$\Gamma_{\sigma\nu}^\mu = \frac{1}{2}g^{\mu\rho} \left(\frac{\partial g_{\sigma\rho}}{\partial x^\nu} + \frac{\partial g_{\nu\rho}}{\partial x^\sigma} - \frac{\partial g_{\sigma\nu}}{\partial x^\rho} \right). \quad (3.13)$$

With this definition Newton's second law has been made generally covariant.

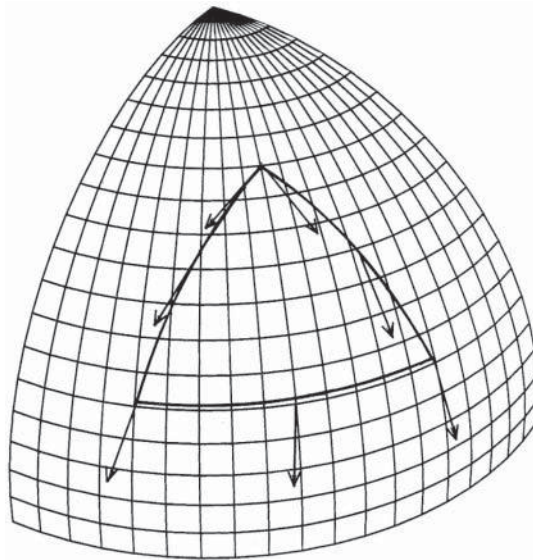


Figure 3.4 Parallel transport of a vector around a closed contour on a curved surface.

The path of a test body in free fall follows from Equation (3.12) by requiring that no forces act on the body, $F^\mu = 0$. Making the replacement

$$P^\mu = m \frac{dx^\mu}{d\tau},$$

the relativistic equation of motion of the test body, its *geodesic equation*, can be written

$$\frac{d^2 x^\mu}{d\tau^2} + \Gamma_{\sigma\nu}^\mu \frac{dx^\sigma}{d\tau} \frac{dx^\nu}{d\tau} = 0. \quad (3.14)$$

In an inertial frame the metric is flat, the metric tensor is a constant everywhere, $g_{\mu\nu}(x) = \eta_{\mu\nu}$, and thus the space-time derivatives of the metric tensor vanish:

$$\frac{\partial g_{\mu\nu}(x)}{\partial x^\rho} = 0. \quad (3.15)$$

It then follows from Equation (3.13) that the affine connections also vanish, and the covariant derivatives equal the simple space-time derivatives.

Going from an inertial frame at x to an accelerated frame at $x + \Delta x$ the expressions for $g_{\mu\nu}(x)$ and its derivatives at x can be obtained as the Taylor expansions

$$g_{\mu\nu}(x + \Delta x) = \eta_{\mu\nu} + \frac{1}{2} \frac{\partial^2 g_{\mu\nu}(x)}{\partial x^\rho \partial x^\sigma} \Delta x^\rho \Delta x^\sigma + \dots$$

and

$$\frac{\partial g_{\mu\nu}(x + \Delta x)}{\partial x^\rho} = \frac{\partial^2 g_{\mu\nu}(x)}{\partial x^\rho \partial x^\sigma} \Delta x^\sigma + \dots$$

The description of a curved space-time thus involves second derivatives of $g_{\mu\nu}$, at least. (Only in a very strongly curved space-time would higher derivatives be needed.)

Recall the definition of the noncovariant Gaussian curvature K in Equation (2.31) defined on a curved two-dimensional surface. In a higher-dimensional space-time, curvature has to be defined in terms of more than just one parameter K . It turns out that curvature is most conveniently defined in terms of the fourth-rank *Riemann tensor*

$$R_{\beta\gamma\sigma}^\alpha = \frac{\partial \Gamma_{\beta\sigma}^\alpha}{\partial x^\gamma} - \frac{\partial \Gamma_{\beta\gamma}^\alpha}{\partial x^\sigma} + \Gamma_{\rho\gamma}^\alpha \Gamma_{\beta\sigma}^\rho - \Gamma_{\rho\sigma}^\alpha \Gamma_{\beta\gamma}^\rho. \quad (3.16)$$

In four-space this tensor has 256 components, but most of them vanish or are not independent because of several symmetries and antisymmetries in the indices. Moreover, an observer at rest in the comoving Robertson–Walker frame will only need to refer to spatial curvature. In a spatial n -manifold, $R_{\beta\gamma\delta}^\alpha$ has only $n^2(n^2 - 1)/12$ nonvanishing components, thus six in the spatial three-space of the Robertson–Walker metric. On the two-sphere there is only one component, which is essentially the Gaussian curvature K .

Another important tool related to curvature is the second rank *Ricci tensor* $R_{\beta\gamma}$, obtained from the Riemann tensor by a summing operation over repeated indices, called *contraction*:

$$R_{\beta\gamma} = R_{\beta\gamma\alpha}^\alpha = \delta_\alpha^\sigma R_{\beta\gamma\sigma}^\alpha = g^{\alpha\sigma} R_{\beta\gamma\sigma}^\alpha. \quad (3.17)$$

This n^2 -component tensor is symmetric in the two indices, so it has only $\frac{1}{2}n(n+1)$ independent components. In four-space the ten components of the Ricci tensor lead

to Einstein's system of ten gravitational equations as we shall see later. Finally, we may sum over the two indices of the Ricci tensor to obtain the *Ricci scalar* R :

$$R = g^{\beta\gamma} R_{\beta\gamma}, \quad (3.18)$$

which we will need later. Actually we first met contraction in its simplest form, the scalar product of two vectors [Equation (3.9)].

3.3 The Einstein Equation

Realizing that the space we live in was not flat, except locally and approximately, Einstein proceeded to combine the equivalence principle with the requirement of general covariance. The inhomogeneous gravitational field near a massive body being *equivalent* to (a patchwork of flat frames describing) a curved space-time, the laws of nature (such as the law of gravitation) have to be described by *generally covariant* tensor equations. Thus the law of gravitation has to be a covariant relation between mass density and curvature. Einstein searched for the simplest form such an equation may take.

In the analysis of classical fields as carriers of forces one describes the dynamics by equations of motion. A convenient way to derive the equations of motion is achieved by introducing the *Lagrangian* L which is the difference between kinetic and potential energies, and the *action* $\int L dt$. Then the *action principle* is an extremal of the action,

$$\delta \int L dt = 0, \quad (3.19)$$

which delivers the equations of motion. This is also called the *principle of least action*.

But here we are interested in deriving the Einstein equation of gravitation. We start with a Lagrangian defined in terms of a *Lagrangian density* L , and an action of the relativistic form

$$\delta \int L d^4x^\mu = 0. \quad (3.20)$$

The Einstein–Hilbert Action. In general relativity, the action is usually assumed to be a functional of the metric tensor $g_{\mu\nu}$ and the affine connexions (3.13). This is the *Einstein–Hilbert action* proposed in 1915 by *David Hilbert* (1862–1943),

$$S = \int \left[\frac{1}{2\kappa} R + \mathcal{L}_M \right] \sqrt{-g} d^4x, \quad (3.21)$$

where $\kappa = 8\pi G/c^4$, R is the Ricci scalar 3.18 and $g = \det(g_{\mu\nu})$. The term \mathcal{L}_M describes any matter fields appearing in the theory.

Note that the integral S is defined over all of spacetime d^4x , which is of course a simplification. General relativity assumes the Copernican principle to be true, but this is only accurate on small scales and modifications are needed at larger scales. The scale at which the assumption breaks down is still debated but unknown.

The action principle then requires the variation of this action with respect to the inverse metric $g^{\mu\nu}$ to vanish, $\delta\mathcal{L} = 0$. Spelled out,

$$\int \left[\frac{1}{2\kappa} \frac{\delta\sqrt{-g}R}{\delta g^{\mu\nu}} + \frac{\delta(\sqrt{-g}\mathcal{L}_M)}{\delta g^{\mu\nu}} \right] \delta g^{\mu\nu} d^4x. \quad (3.22)$$

Multiplying the integrand by $2\kappa/\sqrt{-g}$ the square bracket is a function that can be set equal to zero for any variations $\delta g^{\mu\nu}$, and we are left with an equation with three terms for the motion of the metric field:

$$\frac{\delta R}{\delta g^{\mu\nu}} + \frac{R}{\sqrt{-g}} \frac{\delta\sqrt{-g}}{\delta g^{\mu\nu}} = -\frac{2\kappa}{\sqrt{-g}} \frac{\delta(\sqrt{-g}\mathcal{L}_M)}{\delta g^{\mu\nu}}. \quad (3.23)$$

The right hand side of this equation is a tensor which can be chosen to be proportional the stress-energy tensor $T_{\mu\nu}$ since it contains all the matter terms \mathcal{L}_M .

The first term on the left is the variation of the Ricci scalar. From its definition in Equation (3.18) we have

$$\delta R = R_{\mu\nu}\delta g^{\mu\nu} + g^{\mu\nu}\delta R_{\mu\nu}. \quad (3.24)$$

One can show that the second term does not contribute when integrated over the whole space-time so we have the result

$$\frac{\delta R}{\delta g^{\mu\nu}} = R_{\mu\nu}. \quad (3.25)$$

The second term in Equation (3.23) requires the rule for the variation of a determinant,

$$\delta g = g g^{\mu\nu}\delta g_{\mu\nu}. \quad (3.26)$$

Using this we get

$$\delta\sqrt{-g} = -\frac{\delta g}{2\sqrt{-g}} = -\frac{1}{2}\sqrt{-g}(g_{\mu\nu}\delta g^{\mu\nu}). \quad (3.27)$$

It follows that the left-side terms in Equation (3.22) are

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} \quad (3.28)$$

Equation (3.28) expresses that the energy densities, pressures and shears embodied by the stress-energy tensor determine the geometry of space-time, which, in turn, determines the motion of matter. Thus we arrive at the covariant formula for *Einstein gravity*:

$$G_{\mu\nu} = \frac{8\pi G}{c^4}T_{\mu\nu}. \quad (3.29)$$

Stress–Energy Tensor. Let us now turn to the distribution of matter in the Universe. Suppose that matter on some scale can be considered to be continuously distributed as in an *ideal fluid*. The energy density, pressure and shear of a fluid of nonrelativistic

matter are compactly described by the *stress–energy tensor* $T_{\mu\nu}$ with the following components.

1. The time–time component T_{00} is the energy density ρc^2 , which includes the mass as well as internal and kinetic energies.
2. The diagonal space–space components T_{ii} are the pressure components in the i direction p^i , or the momentum components per unit area.
3. The time–space components cT_{0i} are the energy flow components per unit area in the i direction.
4. The space–time components cT_{i0} are the momentum densities in the i direction.
5. The nondiagonal space–space components T_{ij} are the shear components of the pressure p^i in the j direction.

It is important to note that the stress–energy tensor is of rank 2 and is symmetric, thus it has ten independent components in four-space. However, a comoving observer in the Robertson–Walker space-time following the motion of the fluid sees no time–space or space–time components. Moreover, we can invoke the cosmological principle to neglect the anisotropic nondiagonal space–space components. Thus the stress–energy tensor can be cast into purely diagonal form:

$$T_{\mu\mu} = (p + \rho c^2) \frac{v_\mu v_\mu}{c^2} - p g_{\mu\mu}. \quad (3.30)$$

In particular, the time–time component T_{00} is ρc^2 . The conservation of energy and three-momentum, or equivalently the conservation of four-momentum, can be written

$$\frac{DT_{\mu\nu}}{Dx_\nu} = 0. \quad (3.31)$$

Thus the stress–energy tensor is divergence free.

Taking $T_{\mu\nu}$ to describe relativistic matter, one has to pay attention to its Lorentz transformation properties, which differ from the classical case. Under Lorentz transformations the different components of a tensor do not remain unchanged, but become dependent on each other. Thus the physics embodied by $T_{\mu\nu}$ also differs: the gravitational field does not depend on mass densities alone, but also on pressure. All the components of $T_{\mu\nu}$ are therefore responsible for warping the space-time.

The stress-energy tensor $T_{\mu\nu}$ is the sum of the stress-energy tensors for the various components of energy, baryons, radiation, neutrinos, dark matter and possible other forms. Einstein’s formula [Equation (3.29)] expresses that the energy densities, pressures and shears embodied by the stress-energy tensor determine the geometry of space-time, which, in turn, determines the motion of matter.

Energy of Gravity Waves. The electromagnetic field which is described by Maxwell’s field equation has a source, the electric charge. In contrast, the gravitational field $g_{\mu\nu}$

does not have a source, it is its own source. The only source of energy in the Einstein Equation (3.26) is matter in the stress-energy tensor. In the absence of matter $G_{\mu\nu} = 0$. But the question then arises how to include the energy of the gravitational field itself. Since gravitational waves are expected to cause physical effects at some distance from the source, the Einstein equation is clearly lacking something.

This lack arises because the Einstein equation describes a linearized theory, it is formulated using the Ricci tensor which is also linearized and which vanishes in the absence of a stress-energy tensor. The remedy is to introduce higher-order corrections in Equation (3.26) [4]. One replaces the gravity field $g_{\mu\nu}$ by $g_{\mu\nu}(I) + g_{\mu\nu}(II)$ where (I) and (II) refer to first and second order terms, and where the background metric $g_{\mu\nu}(I)$ is not regarded as static, but as responding to the gravity waves. Thus in the absence of matter,

$$G_{\mu\nu} = G_{\mu\nu}(I) + G_{\mu\nu}(II) = 0, \quad (3.32)$$

one has $G_{\mu\nu}(I) = -G_{\mu\nu}(II)$. At large scale where one can neglect the wavelengths of the gravity waves the first order description is adequate, at shorter scales the influence of gravity waves is described by $G_{\mu\nu}(II)$ and higher order terms.

Carrying out this argument in full detail leads to replacing the Einstein tensor by a pseudotensor [4].

3.4 Weak Field Limit

The starting point is Newton's law of gravitation, because this has to be true anyway in the limit of very weak fields. From Equation (1.27), the gravitational force experienced by a unit mass at distance r from a body of mass M and density ρ is a vector in three-space

$$\ddot{\mathbf{r}} = \mathbf{F} = -\frac{GM\mathbf{r}}{r^3},$$

in component form ($i = 1, 2, 3$)

$$\frac{d^2 x^i}{dt^2} = F^i = -\frac{GMx^i}{r^3}. \quad (3.33)$$

Let us define a scalar *gravitational potential* ϕ by

$$\frac{\partial \phi}{\partial x^i} = -F^i.$$

This can be written more compactly as

$$\nabla \phi = -\mathbf{F}. \quad (3.34)$$

Integrating the flux of the force \mathbf{F} through a spherical surface surrounding M and using Stokes's theorem, one can show that the potential ϕ obeys Poisson's equation

$$\nabla^2 \phi = 4\pi G\rho. \quad (3.35)$$

Let us next turn to the relativistic equation of motion (3.14). In the limit of weak and slowly varying fields for which all time derivatives of $g_{\mu\nu}$ vanish and the

(spatial) velocity components $dx^i/d\tau$ are negligible compared with $dx^0/d\tau = c dt/d\tau$, Equation (3.14) reduces to

$$\frac{d^2x^\mu}{d\tau^2} + c^2\Gamma_{00}^\mu\left(\frac{dt}{d\tau}\right)^2 = 0. \quad (3.36)$$

From Equation (3.13) these components of the affine connection are

$$\Gamma_{00}^\mu = -\frac{1}{2}g^{\mu\rho}\frac{\partial g_{00}}{\partial x^\rho},$$

where g_{00} is the time–time component of $g_{\mu\nu}$ and the sum over ρ is implied.

In a weak static field the metric is almost that of flat space-time, so we can approximate $g_{\mu\nu}$ by

$$g_{\mu\nu} = \eta_{\mu\nu} + h_{\mu\nu},$$

where $h_{\mu\nu}$ is a small increment to $\eta_{\mu\nu}$. To lowest order in $h_{\mu\nu}$ we can then write

$$\Gamma_{00}^\mu = -\frac{1}{2}\eta^{\mu\rho}\frac{\partial h_{00}}{\partial x^\rho}. \quad (3.37)$$

Inserting this expression into Equation (3.36), the equations of motion become

$$\frac{d^2\mathbf{x}}{d\tau^2} = -\frac{1}{2}\left(\frac{dt}{d\tau}\right)^2 c^2\nabla h_{00}, \quad (3.38)$$

$$\frac{d^2t}{d\tau^2} = 0. \quad (3.39)$$

Dividing Equation (3.38) by $(dt/d\tau)^2$ we obtain

$$\frac{d^2\mathbf{x}}{dt^2} = -\frac{1}{2}c^2\nabla h_{00}. \quad (3.40)$$

Comparing this with the Newtonian equation of motion (3.31) in the x^i direction we obtain the value of the time–time component of $h_{\mu\nu}$,

$$h_{00} = 2\frac{\phi}{c^2},$$

from which it follows that

$$g_{00} = 1 + 2\frac{\phi}{c^2} = 1 - \frac{2GM}{c^2r}. \quad (3.41)$$

We can now put several things together: replacing ρ in the field equation (3.35) T_{00}/c^2 and substituting ϕ from Equation (3.41) we obtain a field equation for weak static fields generated by nonrelativistic matter:

$$\nabla^2 g_{00} = \frac{8\pi G}{c^4}T_{00}. \quad (3.42)$$

Let us now assume with Einstein that the right-hand side could describe the source term of a relativistic field equation of gravitation if we made it generally covariant. This suggests replacing T_{00} with $T_{\mu\nu}$. In a matter-dominated universe where the gravitational field is produced by massive stars, and where the pressure between stars is negligible, the only component of importance is then T_{00} .

The left-hand side of Equation (3.42) is not covariant, but it does contain second derivatives of the metric, albeit of only one component. Thus it is already related to curvature. The next step would be to replace $\nabla^2 g_{00}$ by a tensor matching the properties of $T_{\mu\nu}$ on the right-hand side.

- (i) It should be of rank two.
- (ii) It should be related to the Riemann curvature tensor $R_{\alpha\beta\gamma\sigma}$. We have already found a candidate in the Ricci tensor $R_{\mu\nu}$ in Equation (3.17).
- (iii) It should be symmetric in the two indices. This is true for the Ricci tensor.
- (iv) It should be divergence-free in the sense of covariant differentiation. This is not true for the Ricci tensor, but a divergence-free combination can be formed with the Ricci scalar R in Equation (3.18).

The Einstein tensor $G_{\mu\nu}$ contains only terms which are either quadratic in the first derivatives of the metric tensor or linear in the second derivatives.

For weak stationary fields produced by nonrelativistic matter, G_{00} indeed reduces to $\nabla^2 g_{00}$. The Einstein tensor vanishes for flat space-time and in the absence of matter and pressure, as it should. Thus the problems encountered by Newtonian mechanics and discussed at the end of Section 1.7 have been resolved in Einstein's theory. The recession velocities of distant galaxies do not exceed the speed of light, and effects of gravitational potentials are not felt instantly, because the theory is relativistic. The discontinuity of homogeneity and isotropy at the boundary of the Newtonian universe has also disappeared because four-space is unbounded, and because space-time in general relativity is generated by matter and pressure. Thus space-time itself ceases to exist where matter does not exist, so there cannot be any boundary between a homogeneous universe and a void outside space-time.

Problems

1. Derive the Taylor expansions quoted below Equation (3.15).
2. Derive Newton's second law in the generally covariant form Equation (3.12).
3. Show that the Ricci tensor $R_{\beta\gamma}$ is symmetric.
4. Show that $G_{\mu\nu}$ in Equation (3.28) is divergence-free in the sense of covariant differentiation.

References

- [1] Kenyon, I. R. 1990 *General relativity*. Oxford University Press, Oxford.
- [2] Shore, G. M. 2002 *Nuclear Phys. B* **633**, 271.
- [3] Peebles, P. J. E. 1993 *Principles of physical cosmology*. Princeton University Press, Princeton, NJ.
- [4] Peacock, J. A. 1999 *Cosmological physics*. Cambridge University Press, Cambridge.

4

Tests of General Relativity

Before we proceed to construct a theory for a possible description of the Universe we should take a look in this chapter at several phenomena serving as tests of general relativity. Recall the statistical meaning of tests, discussed at the end of Section 1.4.

In Section 4.1 we describe the classical tests of general relativity which provided convincing evidence early on that the theory was valid, if only in the weak field limit. In Section 4.2 we describe the precision measurements of properties of two binary pulsars which testify to the correctness of general relativity. A very important gravitational phenomenon is gravitational lensing, encountered already in the early observations of starlight deflected by the passage near the Sun's limb. The strong lensing of distant galaxies and quasars by interposed galaxy clusters, which is discussed in Section 4.3, has become a tool for studying the internal structure of clusters. and the distribution of dark matter Weak lensing is a tool for studying the large-scale distribution of matter in the Universe, in particular the dark matter which is not observed by any radiation.

The existence of gravitational radiation, already demonstrated in the case of the binary pulsar, is an important prediction of general relativity. However, it remains a great challenge to observe this radiation directly. How to do this will be described in Section 4.4.

The extremely important case of black holes as tests of general relativity is left for the next Chapter, Section 5.4.

4.1 The Classical Tests

The classical testing ground of theories of gravitation, Einstein's among them, is celestial mechanics within the Solar System. Ideally one should consider the full

many-body problem of the Solar System, a task which one can readily characterize as impossible. Already the relativistic two-body problem presents extreme mathematical difficulties. Therefore, all the classical tests of general relativity treated only the one-body problem of the massive Sun influencing its surroundings, and only in the weak field limit.

Note that the expansion of the Universe and Hubble's linear law [Equation (4.15)] are not tests of general relativity. Objects observed at wavelengths ranging from radio to gamma rays are close to isotropically distributed over the sky. Either we are close to a center of spherical symmetry—an anthropocentric view—or the Universe is close to homogeneous. In the latter case, and if the distribution of objects is expanding so as to preserve homogeneity and isotropy (this is local Lorentz invariance), the recession velocities satisfy Hubble's law.

Mercury's Perihelion Shift. The earliest phenomenon requiring general relativity for its explanation was noted in 1859, 20 years before Einstein's birth. The French astronomer *Urban Le Verrier* (1811–1877) found that something was wrong with the planet Mercury's elongated elliptical orbit. As the innermost planet it feels the solar gravitation very strongly, but the orbit is also perturbed by the other planets. The total effect is that the elliptical orbit is nonstationary: it precesses slowly around the Sun. The locus of Mercury's orbit nearest the Sun, the *perihelion*, advances 574" (seconds of arc) per century. This is calculable using Newtonian mechanics and Newtonian gravity, but the result is only 531" which is 43" too little. Le Verrier, who had already successfully predicted the existence of Neptune from perturbations in the orbit of Uranus, suspected that the discrepancy was caused by a small undetected planet inside Mercury's orbit, which he named Vulcan. That prediction was, however, never confirmed.

With the advent of general relativity the calculations could be remade. This time the discrepant 43" were successfully explained by the new theory, which thereby gained credibility. This counts as the first one of three 'classical' tests of general relativity. For details on this test as well as on most of the subsequent tests see, for example, [1] and [2].

Also, the precessions of Venus and Earth have been put to similar use, and within the Solar System many more consistency tests have been done, based on measurements of distances and other orbital parameters.

Deflection of Star Light. The second classical test was the predicted deflection of a ray of light passing near the Sun. A consequence of the relativistic phenomenon of light rays bending around gravitating masses is that masses can serve as *gravitational lenses* if the distances are right and the gravitational potential is sufficient. Newton discussed the possibility that celestial bodies could deflect light (in 1704), and the astronomer Soldner published a paper (in 1804) in which he obtained the correct Newtonian deflection angle by the Sun, assuming that light was corpuscular in nature. Einstein published the general relativistic calculation of this deflection only in 1936, and it was not until 1979 that a suitable solar eclipse occurred which permitted astronomers to see the effect.

We shall return to a fuller presentation of gravitational lensing in Section 4.3.

Timekeeping in Gravitational Fields. The third classical test was the gravitational shift of atomic spectra, first observed by *John Evershed* in 1927. The frequency of emitted radiation makes atoms into clocks. In a strong gravitational field these clocks run slower, so the atomic spectra shift towards lower frequencies. This is an effect which we already met in Equation (3.1): light emerging from a star with mass M is gravitationally redshifted in proportion to M . Evershed observed the line shifts in a cloud of plasma ejected by the Sun to an elevation of about 72 000 km above the photosphere and found an effect only slightly larger than that predicted by general relativity. Modern observations of atoms radiating above the photosphere of the Sun have improved on this result, finding agreement with theory at the level of about 2.1×10^{-6} . Similar measurements have been made in the vicinity of more massive stars such as Sirius.

Thus time passes faster at higher elevations above the ground so that the biological time of an astronaut in a space craft passes faster than on the ground—the famous *twin paradox*.

The *Global Positioning System* was originally proposed as a test of general relativity using accurate atomic clocks in orbit inside space satellites. Position determinations with radio signals from GPS satellites in well-known orbits are based on the Doppler effect which is not relativistic. However, calculations based on general relativity showed that the clocks in the satellites would be seen by terrestrial observers to run 38 microseconds faster per day, which had to be corrected for.

The most accurate timekeeping devices built today are atomic clocks with a precision of the order of 10^{-18} s and can measure a vertical separation of 33 cm.

Radio Signal Delay. Many experiments have studied the effects of changes in a gravitational potential on the rate of a clock or on the frequency of an electromagnetic signal. The so-called ‘fourth’ test of general relativity, which was conceived by *I. I. Shapiro* in 1964 and carried out successfully in 1971 and later, deserves a special mention. This is based on the prediction that an electromagnetic wave suffers a time delay when traversing an increased gravitational potential.

The fourth test was carried out with the radio telescopes at the Haystack and Arecibo observatories by emitting radar signals towards Mercury, Mars and, notably, Venus, through the gravitational potential of the Sun. The round-trip time delay of the reflected signal was compared with theoretical calculations. Further refinement was achieved later by posing the Viking Lander on the Martian surface and having it participate in the experiment by receiving and retransmitting the radio signal from Earth. This experiment found the ratio of the delay observed to the delay predicted by general relativity to be 1.000 ± 0.002 .

4.2 Binary Pulsars

The most important tests have been carried out on the radio observations of pulsars that are members of binary pairs of two neutron stars or one neutron star spinning around a white dwarf. The PSR 1913 + 16 discovered in 1974 by *R. A. Hulse* and

J. H. Taylor, for which they received the Nobel Prize in 1993, is a neutron star–neutron star pair. Pulsars are rapidly rotating (up to 700 times per second), strongly magnetized neutron stars with a surface gravity 10^{11} stronger than the Earth's and with magnetic fields ranging from 10^{11} to 10^{15} stronger than the Earth's.

If the magnetic dipole axis does not coincide with the axis of rotation (just as is the case with Earth), the star would radiate copious amounts of energy along the magnetic dipole axis. These beams at radio frequencies precess around the axis of rotation like the searchlights of a beacon. As the beam sweeps past our line of sight, it is observable as a pulse with the period of the rotation of the star. Hulse, Taylor and collaborators at Arecibo have demonstrated that pulsars are the extremely stable clocks, the time variation of the PSR 1913 + 16 is about 10^{-14} on timescales of 6–12 months. The reason for this stability is the intense self-gravity of a neutron star, which makes it almost undeformable until, in a binary pair, the very last few orbits when the pair coalesce into one star.

The neutron stars in the binary system PSR 1913 + 16 rotate, in addition to their individual spins, also around their common center of mass in a quite eccentric orbit. One of the binary stars is a pulsar, sweeping in our direction with a period of 59 ms, and the binary period of the system is determined to be 7.751 939 337 h. The radial velocity curve as a function of time is known, and from this one can deduce the masses m_1 and m_2 of the binary stars to a precision of 0.000 5, as well as the parameters of a Keplerian orbit: the eccentricity and the semi-major axis.

Pulsars lose energy due to its emission of a relativistic wind and electromagnetic radiation. But the binary system does not behave exactly as expected in Newtonian astronomy, hence the deviations provide several independent confirmations of general relativity. The largest relativistic effect is the apsidal motion of the orbit, which is analogous to the advance of the perihelion of Mercury. A second effect is the counterpart of the relativistic clock correction for an Earth clock. The light travel time of signals from the pulsar through the gravitational potential of its companion provides a further effect.

During the first 17 years of observations the team observed a steadily accumulating change of orbital phase of the binary system, which must be due to the loss of orbital rotational energy by the emission of *gravitational radiation*. This rate of change can be calculated since one knows the orbital parameters and the star masses so well. The calculations based on Einstein's general relativity agree to within 1% with the measurements. This was the first observation of gravitational radiation, although it is indirect, since we as yet have no detector with which to receive such waves. The result was also an important check on competing gravitational theories, several of which were ruled out.

To test general relativity in the strong field regime and to search for signals of gravitational radiation one requires a massive pulsar spinning in a tight orbit around a white dwarf. The best example is the pulsar PSR J0348+0432 discovered in 2011, with a mass of $2.01M_{\odot}$ spinning at 3.9 ms in a 2.46-hr orbit with a low-mass companion, $\approx 2M_{\odot}$. The pulsar is losing energy due to gravitational radiation at a rate in agreement with Einstein's general relativity to within $5 \pm 18\%$.

A more recently discovered millisecond pulsar is PSR J0337+1715, a hierarchical triple system with two other stars. Strong gravitational interactions are apparent and provide precision timing and multi-wavelength observations. The masses of the pulsar and the two white dwarf companions are $1.4378 \pm 0.0013 M_{\odot}$, $0.19751 \pm 0.0015 M_{\odot}$ and $0.4101 \pm 0.0003 M_{\odot}$, respectively. The unexpectedly coplanar and nearly circular orbits indicate a complex and exotic evolutionary past that differs from those of known stellar systems. The gravitational field of the outer white dwarf strongly accelerates the inner binary containing the neutron star, and the system will thus provide an ideal laboratory in which to test the strong equivalence principle of general relativity.

4.3 Gravitational Lensing

Recall from Equation (3.4) and Section 3.1 that the Strong Equivalence Principle (SEP) causes a photon in a gravitational field to move as if it possessed mass. A particle moving with velocity v past a gravitational point potential or a spherically symmetric potential U will experience an acceleration in the transversal direction resulting in a deflection, also predicted by Newtonian dynamics. The deflection angle α can be calculated from the (negative) potential U by taking the line integral of the transversal gravitational acceleration along the photon's path.

Weak Lensing. In the thin-lens approximation the light ray propagates in a straight line, and the deflection occurs discontinuously at the closest distance. The transversal acceleration in the direction y is then

$$\frac{d^2y}{dt^2} = - \left(1 + \frac{v^2}{c^2} \right) \frac{dU}{dy}. \quad (4.1)$$

In Newtonian dynamics the factor in the brackets is just 1, as for velocities $v \ll c$. This is also true if one invokes SEP alone, which accounts only for the distortion of time. However, the full theory of general relativity requires the particle to move along a geodesic in a geometry where space is also distorted by the gravitational field. For photons with velocity c the factor in brackets is then 2, so that the total deflection due to both types of distortion is doubled.

The gravitational distortion can be described as an effective refraction index,

$$n = 1 - \frac{2}{c^2} U > 1, \quad (4.2)$$

so that the speed of light through the gravitational field is reduced to $v = c/n$. Different paths suffer different time delays Δt compared with undistorted paths:

$$\Delta t = \frac{1}{c} \int_{\text{source}}^{\text{observer}} \frac{2}{c^2} dl. \quad (4.3)$$

For a thin lens, deflection through the small bend angle α in Figure 4.1 may be taken to be instantaneous. The angles θ_I and θ_S specify the observed and intrinsic positions

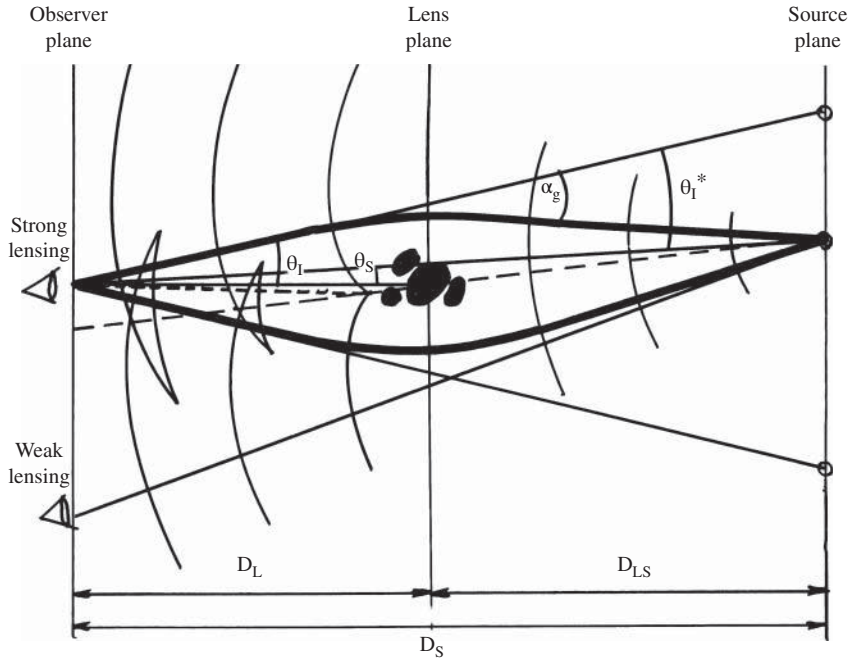


Figure 4.1 A gravitationally lensing cluster. Wavefronts and light rays are shown. The geometry is described in the text.

of the source on the sky, respectively. Weak lensing is defined as events where the potential is weak and the light ray clearly avoids the lensing object (the lower light ray in Figure 4.1). The field equations of GR can then be linearized.

For the special case of a spherically or circularly symmetric gravitating body such as the Sun with mass M inside a radius b , photons passing at distance b of closest approach would be deflected by the angle (Problem 2)

$$\alpha = \frac{4GM}{bc^2}. \quad (4.4)$$

For light just grazing the Sun's limb ($b = 6.96 \times 10^8$ m), the relativistic deflection is $\alpha = 1.750''$, whereas the nonrelativistic deflection would be precisely half of this.

To observe the deflection one needs stars visible near the Sun, so two conditions must be fulfilled. The Sun must be fully eclipsed by the Moon to shut out its intense direct light, and the stars must be very bright to be visible through the solar corona. Soon after the publication of Einstein's theory in 1917 it was realized that such a fortuitous occasion to test the theory would occur on 29 May 1919. The Royal Astronomical Society then sent out two expeditions to the path of the eclipse, one to Sobral in North Brazil and the other one, which included *Arthur S. Eddington* (1882–1944), to the Isle of Principe in the Gulf of Guinea, West Africa.

Both expeditions successfully observed several stars at various distances from the eclipsed Sun, and the angles of deflection (reduced to the edge of the Sun) were

$1.98'' \pm 0.12''$ at Sobral and $1.61'' \pm 0.30''$ at Principe. This confirmed the predicted value of $1.750''$ with reasonable confidence and excluded the Newtonian value of $0.875''$. The measurements have been repeated many times since then during later solar eclipses, with superior results confirming the general relativistic prediction.

The case of starlight passing near the Sun is a special case of a general lensing event, shown in Figure 4.1. For the Sun the distance from lens to observer is small so that the angular size distances are $D_{LS} \approx D_S$, which implies that the actual deflection equals the observed deflection and $\alpha = \theta_I - \theta_S$. In the general case, simple geometry gives the relation between the deflection and the observed displacement as

$$\alpha = \frac{D_S}{D_{LS}}(\theta_I - \theta_S), \quad (4.5)$$

For a lens composed of an ensemble of point masses, the deflection angle is, in this approximation, the vectorial sum of the deflections of the individual point lenses. When the light bending can be taken to be occurring instantaneously (over a short distance relative to D_{LS} and D_L), we have a geometrically thin lens, as assumed in Figure 4.1. Thick lenses are considerably more complicated to analyze.

Strong Lensing. The terms *weak lensing* and *strong lensing* are not defined very precisely. In weak lensing the deflection angles are small and it is relatively easy to determine the true positions of the lensed objects in the source plane from their displaced positions in the observer plane. Strong lensing implies deflection through larger angles by stronger potentials. The images in the observer plane can then become quite complicated because there may be more than one null geodesic connecting source and observer, so that it is not even always possible to find a unique mapping onto the source plane. Strong lensing is a tool for testing the distribution of mass in the lens rather than purely a tool for testing general relativity.

If a strongly lensing object can be treated as a point mass and is positioned exactly on the straight line joining the observer and a spherical or pointlike lensed object, the lens focuses perfectly and the lensed image is a ring seen around the lens, called an *Einstein ring*. The angular size can be calculated by setting the two expressions for α , Equations (4.4) and (4.5), equal, noting that $\theta_S = 0$ and solving for θ_I :

$$\theta_I = \sqrt{\frac{4GM D_{LS}}{c^2 D_L D_S}}. \quad (4.6)$$

For small M the image is just pointlike. In general, the lenses are galaxy clusters or (more rarely) single galaxies that are not spherical and the geometry is not simple, so that the Einstein ring breaks up into an odd number of sections of arc. Each arc is a complete but distorted picture of the lensed object.

In general, the solution of the lensing equation and the formation of multiple images can be found by jointly solving Equations (4.4) and (4.5). Equation (4.4) gives the bend angle $\alpha_g(M_b)$ as a function of the gravitational potential for a (symmetric) mass M_b within a sphere of radius b , or the mass seen in projection within a circle of radius b . From Figure 4.1 we can see that $b = \theta_I \times D_{LS}$, so inserting this into

Equation (4.4) we have

$$\alpha_g(M_b, \theta_l) = \frac{4GM_b}{c^2\theta_l D_{LS}}. \quad (4.7)$$

Equation (4.5) is the fundamental lensing equation giving the geometrical relation between the bend angle α_1 and the source and image positions:

$$\alpha_1(\theta_s, \theta_l) = \frac{D_S}{D_{LS}}(\theta_l - \theta_s). \quad (4.8)$$

There will be an image at an angle θ_l^* that simultaneously solves both equations:

$$\alpha_g(M_b, \theta_l^*) = \alpha_1(\theta_s, \theta_l^*). \quad (4.9)$$

For the case of a symmetric (or point-mass) lens, θ_l^* will be the two solutions to the quadratic

$$2\theta_l^* = \theta_s + \sqrt{\theta_s^2 + \frac{16GM_b D_{LS}}{c^2 D_L D_S}}. \quad (4.10)$$

This reduces to the radius of the Einstein ring when $\theta_s = 0$. The angle corresponding to the radius of the Einstein ring we denote θ_E .

Equation (4.10) describes a pair of hyperbolas so there will always be two images for a point-mass lens. When the source displacement is zero ($\theta_s = 0$) the images will be at the positive and negative roots of Equation (4.6)—the Einstein ring. When θ_s is large the positive root will be approximately equal to θ_s , while the negative root will be close to zero (on the line of sight of the lens). This implies that every point-mass lens should have images of every source, no matter what the separation in the sky. Clearly this is not the case. The reason is that the assumption of a point mass and hyperbolic α_g cannot be maintained for small θ_l .

A more realistic assumption for the mass distribution of a galaxy would be that the density is spherically symmetric, with density as a function of distance from the galactic core, R , given by

$$\rho(R) = \rho_{\text{core}} \left(1 + \frac{R^2}{R_{\text{core}}^2} \right)^{-1}, \quad (4.11)$$

The density is approximately constant (equal to ρ_{core}) for small radii ($R \ll R_{\text{core}}$) and falls off as R^{-2} for large radii. This roughly matches observed mass-density distributions (including dark matter) as inferred from galaxy rotational-velocity observations. The mass will grow like R^3 for $R \ll R_{\text{core}}$ and like R for $R \gg R_{\text{core}}$.

For a nonsymmetric mass distribution, the function α_g can become quite complicated (see, e.g., [4]). Clearly, the problem quickly becomes complex. An example is shown in Figure 4.1, where each light ray from a lensed object propagates as a spherical wavefront. Bending around the lens then brings these wavefronts into positions of interference and self-interaction, causing the observer to see multiple images. The size and shape of the images are therefore changed. From Figure 4.1 one understands how the time delay of pairs of images arises: this is just the time elapsed between

different sheets of the same wavefront. In principle, the time delays in Equation (4.3) provides a tool for measuring H_0 .

Surface Brightness and Microlensing. Since photons are neither emitted nor absorbed in the process of gravitational light deflection, the surface brightness of lensed sources remains unchanged. Changing the size of the cross-section of a light bundle therefore only changes the flux observed from a source and magnifies it at fixed surface-brightness level. For a large fraction of distant quasars the magnification is estimated to be a factor of ten or more. This enables objects of fainter intrinsic magnitudes to be seen. However, lensing effects are very model dependent, so to learn the true magnification effect one needs very detailed information on the structure of the lens.

If the mass of the lensing object is very small, one will merely observe a magnification of the brightness of the lensed object. This is called *microlensing*, and it has been used to search for nonluminous objects in the halo of our Galaxy. One keeps watch over several million stars in the Large Magellanic Cloud (LMC) and records variations in brightness. Some stars are Cepheids, which have an intrinsic variability, so they have to be discarded. A star which is small enough not to emit visible light and which is moving in the halo is expected to cross the diameter of a star in the LMC in a time span ranging from a few days to a couple of months. The total light amplification for all images from a point-mass lens and point source is (Problem 3)

$$A = \frac{1 + \frac{1}{2}x^2}{x\sqrt{1 + \frac{1}{4}x^2}}, \quad x = \frac{\theta_S}{\theta_E}. \quad (4.12)$$

As the relative positions of the source, lens and observer change, θ_S will change. Simple geometrical arguments give θ_S as a function of the relative velocities, and thus the amplification as a function of time (see [4], pp. 106, 118–120). During the occultation, the image of the star behind increases in intensity according to this function and subsequently decreases along the time-symmetric curve. A further requirement is that observations in different colors should give the same time curve. Several such microlensing events have been found in the direction of the LMC, and several hundred in the direction of the bulge of our Galaxy. The number of such occurrences is, however, too small to play any role as nonluminous dark matter later.

Cosmic Shear. The large-scale distribution of matter in the Universe is inhomogeneous in every direction, so one can expect that everything we observe is displaced and distorted by weak lensing. Since the tidal gravitational field, and thus the deflection angles, depend neither on the nature of the matter nor on its physical state, light deflection probes the total projected mass distribution. Lensing in infrared light offers the additional advantage of sensing distant background galaxies, since their number density is higher than in the optical. The idea of mapping the matter distribution using the *cosmic shear* field was already proposed (in 1937) by *Fritz Zwicky* (1898–1974), who also proposed looking for lensing by galaxies rather than by stars.

The ray-tracing process mapping a single source into its image can be expressed by the Jacobian matrix between the source-plane coordinates and the observer-plane coordinates:

$$\mathbf{J}(\alpha) = \begin{pmatrix} 1 - \kappa - \gamma_1 & -\gamma_2 \\ -\gamma_2 & 1 - \kappa + \gamma_1 \end{pmatrix}, \quad (4.13)$$

where κ is the *convergence* of the lens and $\gamma = \gamma_1 + i\gamma_2$ is the shear. The matrix $\mathbf{J}(\alpha)$ transforms a circular source into an ellipse with semi-axes stretched by the factor $(1 - \kappa \pm |\gamma|)^{-1}$. The convergence affects the isotropic magnification or the projected mass density divided by the critical density, whereas the shear affects the shape of the image. The magnification is given by

$$\mu = (\det \mathbf{J})^{-1} = [(1 - \kappa)^2 - \gamma^2]^{-1}. \quad (4.14)$$

Clearly, there are locations where μ can become infinite. These points in the source plane are called *caustics* and they lie on the intersections of *critical curves*.

Weak Lensing Surveys. Background galaxies would be ideal tracers of distortions if they were intrinsically circular. Any measured ellipticity would then directly reflect the action of the gravitational tidal field of the interposed lensing matter, and the statistical properties of the cosmic-shear field would reflect the statistical properties of the matter distribution. But many galaxies are actually intrinsically elliptical, and the ellipses are randomly oriented. These intrinsic ellipticities introduce noise into the inference of the tidal field from observed ellipticities.

The sky is covered with a ‘wall paper’ of faint and distant blue galaxies, about 20 000–40 000 on an area of the size of the full moon. This fine-grained pattern of the sky makes statistical weak-lensing studies possible, because it allows the detection of the coherent distortions imprinted by gravitational lensing on the images of the faint-blue-galaxy population. Large collaborations carrying out such surveys have reported statistically significant observations of cosmic shear and thus of the distribution of interposed lensing dark matter. We shall come back to this discussion in the chapter on dark matter.

To test general relativity versus alternative theories of gravitation, the best way is to probe the gravitational potential far away from visible matter, and weak galaxy–galaxy lensing being a good approach to this end because it is accurate on scales where all other methods fail, and it is simple if galaxies are treated as point masses. Alternative theories may predict an isotropic signal where general relativity predicts an azimuthal variation. The current knowledge favors anisotropy and thus general relativity.

4.4 Gravitational Waves

Einstein noted in 1916 that his general relativity predicted the existence of gravitational radiation, but its possible observation is still in the future. As we explained in Section 4.2, the slowdown of binary pulsars is indirect evidence that the system loses its energy by radiating gravitational waves.

When gravitational waves travel through space-time they produce ripples of curvature, an oscillatory stretching and squeezing of space-time analogous to the tidal effect of the Moon on Earth. Any matter they pass through will feel this effect. Thus a detector for gravitational waves is similar to a detector for the Moon's tidal effect, but the waves act on an exceedingly weaker scale.

Gravitational radiation travels with the speed of light and traverses matter unhindered and unaltered. It may be that the carriers are particles, *gravitons*, with spin $J = 2$, but it is hard to understand how that could be verified. Perhaps, if a theory were found combining gravitation and quantum mechanics, the particle nature of gravitational radiation would be more meaningful.

Tensor Field. In contrast to the electromagnetic field, which is a vector field, the gravitational field is a tensor field. The gravitational analogue of electromagnetic dipole radiation cannot produce any effect because of the conservation of momentum: any dipole radiation caused by the acceleration of an astronomical object is automatically cancelled by an equal and opposite change in momentum in nearby objects. Therefore, gravitational radiation is caused only by nonspherical symmetric accelerations of mass, which can be related to the quadrupole moment, and the oscillatory stretch and squeeze produced is then described by two dimensionless wave fields h_+ and h_\times , which are associated with the gravitational wave's two linear polarizations. If h_+ describes the amplitude of polarization with respect to the x - and y -axes in the horizontal plane, h_\times describes the independent amplitude of polarization with respect to the rotated axes $x + y$ and $x - y$ (see Figure 4.2). The relative tidal effect a detector of length L may observe is then a linear combination of the two wave fields

$$\Delta L/L = a_+ h_+(t) + a_\times h_\times(t) \equiv h(t). \quad (4.15)$$

The proper derivation of the quadrupole formula for the energy loss rate through gravitational radiation of an oscillating body and the spatial strain $h(t)$ caused on bodies elsewhere cannot be carried out here, it requires general relativity to be carried out to high orders of covariant derivation. This complication is a benefit, however, because it renders the detection of gravitational radiation an extremely sensitive test of general relativity.

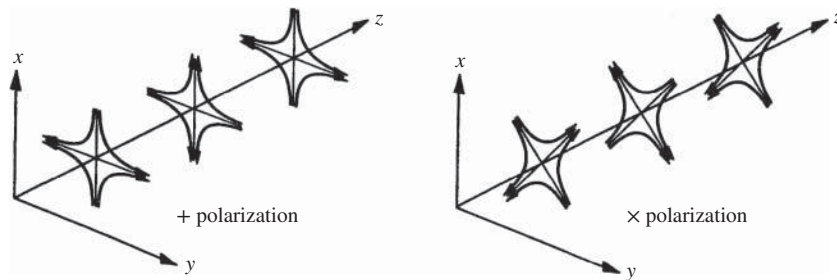


Figure 4.2 The lines of force associated with the two polarizations of a gravitational wave. Reprinted with permission of A. Abramovici *et al.* [3]. Copyright 1992 American Association for the Advancement of Science.

In a Newtonian approximation the strength of the waves from a nonspherical body of mass M , oscillating size $L(t)$, and quadrupole moment $Q(t) \approx ML^2$ at a distance r from Earth is

$$h(t) \approx \frac{G}{c^4 r} \frac{d^2 Q(t)}{dt^2} = \frac{G}{c^4 r} 2Mv(t)^2 = \frac{4G}{c^4 r} E(t), \quad (4.16)$$

where G is the Newtonian constant, v is the internal velocity, and $E = \frac{1}{2} Mv^2$ is the nonspherical part of the internal kinetic energy. The factor c^4 is introduced only to make $h(t)$ dimensionless.

Sources of Gravitational Waves. From this formula one can work out that a nonspherical symmetric supernova collapse at the center of our Galaxy will give rise to waves of amplitude $h \approx 10^{-19}$ causing a subnuclear stretch and squeeze of an object 1 km in length by 10^{-16} m. A spherically symmetric supernova collapse causes no waves. In a catastrophic event such as the collision of two neutron stars or two stellar-mass black holes in which E/c^2 is of the order of one solar mass, Equation (4.16) gives $h \approx 10^{-20}$ at the 16 Mpc distance of the Virgo cluster of galaxies, and $h \approx 10^{-21}$ at a distance of approximately 200 Mpc.

The signals one can expect to observe in the amplitude range $h \approx 10^{-21} - 10^{-20}$ with the present generation of detectors are bursts due to the coalescence of neutron-star binaries during their final minutes and seconds (in the high frequency band $1 - 10^4$ Hz), and periodic waves from slowly merging galactic binaries and extragalactic massive black hole binaries (low-frequency band $10^{-4} - 10^{-2}$ Hz), which are stable over hundreds to millions of years. The timing of millisecond binary pulsars such as the PSR 1913 + 16 belong to the very low-frequency band of $10^{-9} - 10^{-7}$ Hz. In this band, processes in the very early Universe may also act as sources.

Merger waves from superheavy black holes with $10^6 M_\odot$ mass may be so strong that both their direction and their amplitude can be determined by monitoring the waves while the detector rotates around the Sun. This may permit researchers to identify the source with the parallax method and to determine the distance to it with high precision. Combined with redshift measurements of the source, one could determine not only H_0 but even the deceleration parameter q_0 of the Universe. Thus the detection of gravitational waves from black holes would go beyond testing general relativity to determining fundamental cosmological parameters of the Universe.

The dynamics of a hole-hole binary can be divided into three epochs: inspiral, merger and ringdown. The inspiral epoch ends when the holes reach their last stable orbit and begin plunging toward each other. Then the merger epoch commences, during which the binary is a single nonspherical black hole undergoing highly nonlinear space-time oscillations and vibrations of large amplitude. In the ringdown epoch, the oscillations decay due to gravitational wave emission, leaving finally a spherical, spinning black hole.

Gravitational Wave Detection. Detection with huge metal bars as resonant antennas was started by *Joseph Weber* in 1969. These detectors couple to one axis of the

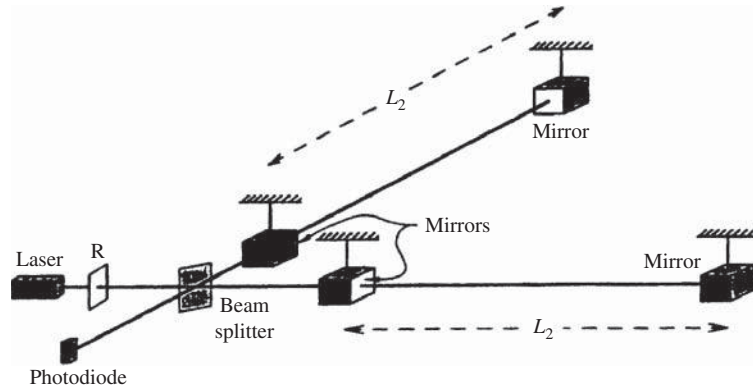


Figure 4.3 A schematic view of a LIGO-type interferometer. Reprinted with permission of A. Abramovici *et al.* [3]. Copyright 1992 American Association for the Advancement of Science.

eigenmodes of the incoming wave, and one then expects to observe a change in the state of oscillation. Today several coordinated and aligned cryogenic bar detectors are in coordinated operation with sensitivities of approximately $10^{-21} \text{ Hz}^{-1/2}$. The detectors are tuned to see approximately 1 ms bursts occurring within a bandwidth of the order of 1 Hz. In order to eliminate random noise, the data from several detectors are analyzed for coincidences.

To improve the signal to noise ratio in the high-frequency range one turns to Michelson interferometers with very long arms. The principle is illustrated in Figure 4.3. A laser beam is split, travels in two orthogonal directions to mirrors, and returns to be recombined and detected. A gravitational wave with either the h_+ or h_x component coinciding with the interferometer axes would lengthen the round-trip distance to one mirror and shorten it to the other. This would be observable as a mismatch of waves upon recombination, and hence as a decrease in the observed combined intensity of the laser. For isolation against mechanical disturbances the optical components are carefully suspended in vacuum. The arm lengths in active terrestrial detectors range from 300 m (TAMA in Japan) and 600 m (GEO600 in Germany) to 3 km (VIRGO in Italy) and 4 km (LIGO at two locations in the United States). Sensitivities of 10^{-21} – $10^{-22} \text{ Hz}^{-1/2}$ can be reached in the high-frequency range. The range is limited to less than approximately 10^4 Hz by photo-electron shot noise in the components of the interferometer.

To study sources in the low-frequency range one has to put the interferometer into space orbiting Earth. This is necessary in order to avoid low-frequency seismic noise on the ground and thermally induced medium-frequency motion in the atmosphere. The spectacular solution is the detector LISA (Laser Interferometer Space Antenna) consisting of three identical spacecraft, forming an equilateral triangle in space, with sidelength 5 million km, trailing Earth by 20° in a heliocentric orbit (see Figure 4.4). From each spacecraft a 1 W beam is sent to the two other remote spacecrafts via a telescope, is reflected by a platinum–gold cubic test mass, and the same telescopes are then used to focus the very weak returning beams. The interference signals

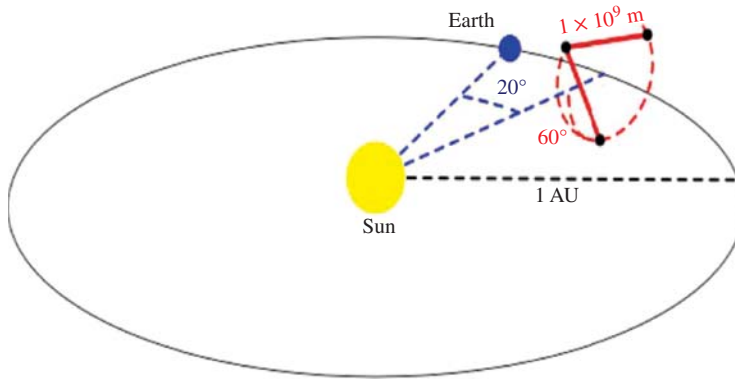


Figure 4.4 eLISA orbits in the Solar System [4]. Reproduced with permission of Pau Amaro-Seoane and the LISA consortium. (See plate section for color version.)

from each arm are combined by on-board computers to perform the multiple-arm interferometry required to cancel the phase-noise common to all three arms. Fluctuations in the optical paths between the test masses can be measured to sub-angstrom precision, which, when combined with the large separation of the spacecraft, allowed LISA to detect gravitational-wave strain down to a level of order 10^{-23} in one year of observation, with a signal to noise ratio of 5.

LISA targets high-priority astronomy such as massive black holes, stellar evolution, the high-redshift universe and cosmology. By 2011 LISA had identified 20 000 individual sources, almost all resolvable binary white dwarfs, but no black hole—black hole mergers and no signals of gravitational radiation [4]. Its follower, the European-led variant named eLISA (see Figure 4.4 [4]) is scheduled to be launched before 2022.

Problems

1. Calculate the gravitational redshift in wavelength for the 769.9 nm potassium line emitted from the Sun's surface [1].
2. Derive the deflection angle [Equation (4.4)] using Equations (4.26) and (4.5).
3. Derive Equation (4.12). What are the amplification of the individual images [4]?
4. Draw the bend angle due to gravitational potential α_g and lens geometry α_l , and zero (dashed line) for a lens of mass $M \approx 7.2 \times 10^{11} M_\odot$, with $D_S \approx 1.64$ Gpc ($z = 3.4$ in a FRW cosmology with $\Omega_m = 0.3$ and $\Omega_\lambda = 0.7$), $D_L \approx 1.67$ Gpc ($z = 0.803$), $D_{LS} \approx 0.96$ Gpc and $\alpha_S \approx 0.13''$. (Note that since distances are angular diameter distances, $D_S \neq D_{LS} + D_L$.) Take the lens to be (a) a point mass and (b) a spherical mass distribution with density given by Equation (4.11) for an ideal galaxy. This roughly corresponds to the parameters associated with the 'Einstein Cross' lensed quasar found by the Hubble telescope, HST 14176 + 5226.

References

- [1] Kenyon, I. R. 1990 *General relativity*. Oxford University Press, Oxford.
- [2] Will, C. M. 1993 *Theory and experiment in gravitational physics*, revised edn. Cambridge University Press, Cambridge.
- [3] Abramovici, A. *et al.* 1992 *Science* **256**, 325.
- [4] Amaro-Seoane, P. *et al.* 2013 *G.W. Notes* **6**, 4–110.

5

Cosmological Models

In Section 5.1 we turn to the ‘concordance’ or Friedmann–Lemaître–Robertson–Walker (FLRW) model of cosmology, really only a paradigm based on Friedmann’s and Lemaître’s equations and the Robertson–Walker metric, which takes both energy density and pressure to be functions of time in a Copernican universe. Among the solutions are the Einstein universe and the Einstein–de Sitter universe, both now known to be wrong, as we shall see in Section 5.4, and the currently accepted Friedmann–Lemaître universe, which includes a positive cosmological constant.

In Section 5.2 we describe the de Sitter model, which does not apply to the Universe at large as we see it now, but which may have dominated the very early universe, and which may be the correct description for the future.

In Section 5.3 we introduce the Schwarzschild solution to the Einstein equation. This takes us to black holes in Section 5.4.

In Section 5.5 we briefly present extensions of general relativity.

5.1 Friedmann–Lemaître Cosmologies

Let us now turn to our main subject, a model describing our homogeneous and isotropic Universe for which the Robertson–Walker metric in Equation (2.32) was derived. Recall that it could be written as a 4×4 tensor with nonvanishing components [Equation (2.33)] on the diagonal only, and that it contained the curvature parameter k .

Friedmann’s Equations. The stress–energy tensor $T_{\mu\nu}$ entering on the right-hand side of the Einstein Equations (3.29) was given by Equation (3.30) in its diagonal form. For a comoving observer with velocity four-vector $v = (c, 0, 0, 0)$, the time–time

component T_{00} and the space–space component T_{11} are then

$$T_{00} = \rho c^2, \quad T_{11} = \frac{\rho a^2}{1 - k\sigma^2}, \quad (5.1)$$

taking g_{00} and g_{11} from Equation (2.33). We will not need T_{22} or T_{33} because they just duplicate the results without adding new dynamical information. In what follows we shall denote mass density by ρ and energy density by ρc^2 . Occasionally, we shall use $\rho_m c^2$ to denote specifically the energy density in all kinds of matter: baryonic, leptonic and unspecified dark matter. Similarly we use $\rho_r c^2$ or ε to specify the energy density in radiation.

On the left-hand side of the Einstein Equations (3.29) we need G_{00} and G_{11} to equate with T_{00} and T_{11} , respectively. We have all the tools to do it: the metric components $g_{\mu\nu}$ are inserted into Equation (3.13) for the affine connection, and subsequently we calculate the components of the Riemann tensor from the expression (3.16) using the metric components and the affine connections. This lets us find the Ricci tensor components that we need, R_{00} and R_{11} , and the Ricci scalar from Equations (3.17) and (3.18), respectively. All this would require several pages to work out (see, e.g., [1, 2]), so I only give the result:

$$G_{00} = 3(ca)^{-2} (\dot{a}^2 + kc^2), \quad (5.2)$$

$$G_{11} = -c^{-2} (2a\ddot{a} + \dot{a}^2 + kc^2) (1 - k\sigma^2)^{-1}. \quad (5.3)$$

Here a is the cosmic scale factor $a(t)$ at time t . Substituting Equations (4.1)–(4.3) into the Einstein Equations (3.29) we obtain two distinct dynamical relations for $a(t)$:

$$\frac{\dot{a}^2 + kc^2}{a^2} = \frac{8\pi G}{3}\rho, \quad (5.4)$$

$$\frac{2\ddot{a}}{a} + \frac{\dot{a}^2 + kc^2}{a^2} = -\frac{8\pi G}{c^2}p. \quad (5.5)$$

These equations were derived in 1922 by Friedmann, seven years before Hubble's discovery, at a time when even Einstein did not believe in his own equations because they did not allow the Universe to be static. Friedmann's equations did not gain general recognition until after his death, when they were confirmed by an independent derivation (in 1927) by Georges Lemaitre (1894–1966). For now they will constitute the tools for our further investigations.

The expansion (or contraction) of the Universe is inherent to Friedmann's equations. Equation (5.4) shows that the rate of expansion, \dot{a} , increases with the mass density ρ in the Universe, and Equation (5.5) shows that it may accelerate. Subtracting Equation (5.4) from Equation (5.5) we obtain

$$\frac{2\ddot{a}}{a} = -\frac{8\pi G}{3c^2}(\rho c^2 + 3p), \quad (5.6)$$

which shows that the acceleration decreases with increasing pressure and energy density, whether mass or radiation energy. Thus it is more appropriate to talk about the *deceleration* of the expansion. Equation 5.6 is also called the *Raychaudhuri equation*.

At our present time t_0 when the mass density is ρ_0 , the cosmic scale is 1, the Hubble parameter is H_0 and the density parameter Ω_0 is given by Equation (1.35), Friedmann's equation (5.4) takes the form

$$\dot{a}_0^2 = \frac{8}{3}\pi G\rho_0 - kc^2 = H_0^2\Omega_0 - kc^2, \quad (5.7)$$

which can be rearranged as

$$kc^2 = H_0^2(\Omega_0 - 1). \quad (5.8)$$

It is interesting to note that this reduces to the Newtonian relation (1.35). Thus the relation between the Robertson–Walker curvature parameter k and the present density parameter Ω_0 emerges: to the k values $+1$, 0 and -1 correspond an overcritical density $\Omega_0 > 1$, a critical density $\Omega_0 = 1$ and an undercritical density $0 < \Omega_0 < 1$, respectively. The spatially flat case with $k = 0$ is called the *Einstein–de Sitter universe*.

General Solution. When we generalized from the present H_0 to the time-dependent Hubble parameter $H(t) = \dot{a}/a$ in Equation (2.47), this also implied that the critical density [Equation (1.31)] and the density parameter [Equation (1.35)] became functions of time:

$$\rho_c(t) = \frac{3}{8\pi G}H^2(t), \quad (5.9)$$

$$\Omega(t) = \rho(t)/\rho_c(t). \quad (5.10)$$

Correspondingly, Equation (5.8) can be generalized to

$$kc^2 = H^2a^2(\Omega - 1). \quad (5.11)$$

If $k \neq 0$, we can eliminate kc^2 between Equations (5.8) and (5.11) to obtain

$$H^2a^2(\Omega - 1) = H_0^2(\Omega_0 - 1), \quad (5.12)$$

which we shall make use of later.

It is straightforward to derive a general expression for the solution of Friedmann's Equation (5.4). Inserting kc^2 from Equation (5.8), and replacing $(8\pi G/3)\rho$ by $\Omega(a)H_0^2$, Equation (5.4) furnishes a solution for $H(a)$:

$$H(a) \equiv \dot{a}/a = H_0\sqrt{(1 - \Omega_0)a^{-2} + \Omega(a)}. \quad (5.13)$$

Here we have left the a dependence of $\Omega(a)$ unspecified. As we shall see later, various types of energy densities with different a dependences contribute.

Equation (5.13) can be used to solve for the *lookback time* $t(z)/t_0$ or $t(a)/t_0$ (normalized to the age t_0) since a photon with redshift z was emitted by writing it as an integral equation:

$$\int_0^{t(a)} dt = \int_1^a \frac{da}{aH(a)}. \quad (5.14)$$

The age of the Universe at a given redshift is then $1 - t(z)/t_0$. We shall specify this in more detail later.

Einstein Universe. Consider now the static universe cherished by Einstein. This is defined by $a(t)$ being constant, $a(t_0) = 1$, so that $\dot{a} = 0$ and $\ddot{a} = 0$ and the age of the Universe is infinite. Equations (5.4) and (5.5) then reduce to

$$kc^2 = \frac{8\pi}{3}G\rho_0 = -\frac{8\pi}{c^2}Gp_0. \quad (5.15)$$

In order that the mass density ρ_0 be positive today, k must be $+1$. Note that this leads to the surprising result that the pressure of matter p_0 becomes negative!

Einstein corrected for this in 1917 by introducing a constant Lorentz-invariant term $\lambda g_{\mu\nu}$ into Equation (3.28), where the *cosmological constant* Ω_λ corresponds to a tiny correction to the geometry of the Universe. Equation (3.28) then becomes

$$G_{\mu\nu} = R_{\mu\nu} - \frac{1}{2}g_{\mu\nu}R - \lambda g_{\mu\nu}. \quad (5.16)$$

In contrast to the first two terms on the right-hand side, the $\lambda g_{\mu\nu}$ term does not vanish in the limit of flat space-time. With this addition, Friedmann's equations take the form

$$\frac{\dot{a}^2 + kc^2}{a^2} - \frac{\lambda}{3} = \frac{8\pi G}{3}\rho, \quad (5.17)$$

$$\frac{2\ddot{a}}{a} + \frac{\dot{a}^2 + kc^2}{a^2} - \lambda = -\frac{8\pi G}{c^2}p. \quad (5.18)$$

A positive value of λ curves space-time so as to counteract the attractive gravitation of matter. Einstein adjusted λ to give a static solution, which is called the *Einstein universe*.

The pressure of matter is certainly very small, otherwise one would observe the galaxies having random motion similar to that of molecules in a gas under pressure. Thus one can set $p = 0$ to a good approximation. In the static case when $a = 1$, $\dot{a}_0 = 0$ and $\ddot{a}_0 = 0$, Equation (5.17) becomes

$$kc^2 - \frac{\lambda}{3} = \frac{8\pi G}{3}\rho_0.$$

It follows from this that in a spatially flat Universe

$$\rho_\lambda = \frac{\lambda}{8\pi G} = -\rho_0. \quad (5.19)$$

But Einstein did not notice that the static solution is unstable: the smallest imbalance between λ and ρ would make \ddot{a} nonzero, causing the Universe to accelerate into expansion or decelerate into contraction. This flaw was only noticed by Eddington in 1930, soon after Hubble's discovery, in 1929, of the expansion that caused Einstein to abandon his belief in a static universe and to withdraw the cosmological constant. This he called 'the greatest blunder of my lifetime'.

The Friedmann–Lemaître Universe. If the physics of the vacuum looks the same to any inertial observer, its contribution to the stress–energy tensor is the same as Einstein's cosmological constant λ , as was noted by Lemaître. The λ term in Equation (5.16) is a correction to the geometrical terms in $G_{\mu\nu}$, but the mathematical content of Equations (5.17) and (5.18) are not changed if the λ terms are moved to the

right-hand side, where they appear as corrections to the stress–energy tensor $T_{\mu\nu}$. Then the physical interpretation is that of an ideal fluid with energy density $\rho_\lambda = \lambda/8\pi G$ and negative pressure $p_\lambda = -\rho_\lambda c^2$. When the cosmological constant is positive, the gravitational effect of this fluid is a cosmic repulsion counteracting the attractive gravitation of matter, whereas a negative λ corresponds to additional attractive gravitation.

The cosmology described by Equations (5.17) and (5.18) with a positive cosmological constant is called the *Friedmann–Lemaître universe* or the Concordance model. Such a universe is now strongly supported by observations of a nonvanishing λ , so the Einstein–de Sitter universe, which has $\lambda = 0$, is a dead end.

In a Friedmann–Lemaître universe the total density parameter is conveniently split into a matter term, a radiation term and a cosmological constant term,

$$\Omega_0 = \Omega_m + \Omega_r + \Omega_\lambda, \quad (5.20)$$

where Ω_r and Ω_λ are defined analogously to Equations (1.35) and (5.10) as

$$\Omega_r = \frac{\rho_r}{\rho_c}, \quad \Omega_\lambda = \frac{\lambda}{8\pi G\rho_c} = \frac{\lambda}{3H_0^2}. \quad (5.21)$$

Ω_m , Ω_r and Ω_λ are important dynamical parameters characterizing the Universe. If there is a remainder $\Omega_k \equiv \Omega_0 - 1 \neq 0$, this is called the *vacuum-energy* term.

Using Equation (5.19) we can find the value of λ corresponding to the attractive gravitation of the present mass density:

$$-\lambda = 8\pi G\rho_0 = 3\Omega_0 H_0^2 \approx 1.3 \times 10^{-52} c^2 \text{ m}^{-2}. \quad (5.22)$$

No quantity in physics this small has ever been known before. It is extremely uncomfortable that λ has to be fine-tuned to a value which differs from zero only in the 52nd decimal place (in units of $c = 1$). It would be much more natural if λ were exactly zero. This situation is one of the enigmas which will remain with us to the end of this book. As we shall see, a repulsive gravitation of this kind may have been of great importance during the first brief moments of the existence of the Universe, and it appears that the present Universe is again dominated by a global repulsion.

Energy-Momentum Conservation. Let us study the solutions of Friedmann’s equations in the general case of nonvanishing pressure p . Differentiating Equation (5.4) with respect to time,

$$\frac{d}{dt}(\dot{a}^2 + kc^2) = \frac{8\pi G}{3} \frac{d}{dt}(\rho a^2),$$

we obtain an equation of second order in the time derivative:

$$2\dot{a}\ddot{a} = \frac{8}{3}\pi G(\dot{\rho}a^2 + 2\rho a\dot{a}). \quad (5.23)$$

Using Equation (5.6) to cancel the second-order time derivative and multiplying through by c^2/a^2 , we obtain a new equation containing only first-order time derivatives:

$$\dot{\rho}c^2 + 3H(\rho c^2 + p) = 0. \quad (5.24)$$

This equation does not contain k and λ , but that is not a consequence of having started from Equations (5.4) and (5.5). If, instead, we had started from Equations (5.17) and (5.18), we would have obtained the same equation.

Note that all terms here have dimension of energy density per time. In other words, Equation (5.24) states that the change of energy density per time is zero, so we can interpret it as the *local energy conservation law*. In a volume element dV , $\rho c^2 dV$ represents the local decrease of gravitating energy due to the expansion, whereas $p dV$ is the work done by the expansion. Energy does not have a global meaning in the curved spacetime of general relativity, whereas work does. If different forms of energy do not transform into one another, each form obeys Equation (5.24) separately. The Einstein equation needs to be extended in some way to be able to serve as a global energy conservation law. There is no “correct” way to do it, but many suggestions. We shall not pursue that search here.

As we have seen, Equation (5.24) follows directly from Friedmann’s equations without any further assumptions. But it can also be derived in another way, perhaps more transparently. Let the total energy content in a comoving volume a^3 be

$$E = (\rho c^2 + p)a^3.$$

The expansion is *adiabatic* if there is no net inflow or outflow of energy so that

$$\frac{dE}{dt} = \frac{d}{dt}[(\rho c^2 + p)a^3] = 0. \quad (5.25)$$

If p does not vary with time, changes in ρ and a compensate and Equation (5.24) immediately follows.

Equation (5.24) can easily be integrated,

$$\int \frac{\dot{\rho}(t)c^2}{\rho(t)c^2 + p(t)} dt = -3 \int \frac{\dot{a}(t)}{a(t)} dt, \quad (5.26)$$

if we know the relation between energy density and pressure—the *equation of state* of the Universe.

Entropy Conservation and the Equation of State. In contrast, the law of *conservation of entropy* S is not implied by Friedmann’s equations, it has to be assumed specifically, as we shall demonstrate in Section 5.2,

$$\dot{S} = 0. \quad (5.27)$$

Then we can make an ansatz for the equation of state: let p be proportional to ρc^2 with some proportionality factor w which is a constant in time,

$$p = w\rho c^2. \quad (5.28)$$

Inserting this ansatz into the integral in Equation (5.26) we find that the relation between energy density and scale is

$$\rho(a) \propto a^{-3(1+w)} = (1+z)^{3(1+w)}. \quad (5.29)$$

Here we use z as well as a because astronomers prefer z since it is an observable. In cosmology, however, it is better to use a for two reasons. Firstly, redshift is a property

of light, but freely propagating light did not exist at times when $z \gtrsim 1080$, so z is then no longer a true observable. Secondly, it is possible to describe the future in terms of $a > 1$, but redshift is then not meaningful.

The value of the proportionality factor w in Equations (5.28) and (5.29) follows from the adiabaticity condition. Leaving the derivation of w for a later discussion, we shall anticipate here its value in three special cases of great importance.

Case I. A *matter-dominated* universe filled with nonrelativistic cold matter in the form of pressureless nonradiating dust for which $p = 0$. From Equation 5.28 then, this corresponds to $w = 0$, and the density evolves according to

$$\rho_m(a) \propto a^{-3} = (1+z)^3. \quad (5.30)$$

It follows that the evolution of the density parameter Ω_m is

$$\Omega_m(a) = \Omega_m \frac{H_0^2}{H^2} a^{-3}.$$

Solving for $H^2 a^2 \Omega$ and inserting it into Equation (5.13), one finds the evolution of the Hubble parameter:

$$H(a) = H_0 a^{-1} \sqrt{1 - \Omega_m + \Omega_m a^{-1}} = H_0 (1+z) \sqrt{1 + \Omega_m z}. \quad (5.31)$$

Case II. A *radiation-dominated* universe filled with an ultra-relativistic hot gas composed of elastically scattering particles of energy density ε . Statistical mechanics then tells us that the equation of state is

$$p_r = \frac{1}{3} \varepsilon = \frac{1}{3} \rho_r c^2. \quad (5.32)$$

This evidently corresponds to $w = \frac{1}{3}$, so that the radiation density evolves according to

$$\rho_r(a) \propto a^{-4} = (1+z)^4. \quad (5.33)$$

Case III. The *vacuum-energy* state corresponds to a flat, static universe ($\ddot{a} = 0$, $\dot{a} = 0$) without dust or radiation, but with a cosmological term. From Equations (5.17) and (5.18) we then obtain

$$p_\lambda = -\rho_\lambda c^2, \quad w = -1. \quad (5.34)$$

Thus the pressure of the vacuum energy is negative, in agreement with the definition in Equation (5.19) of the vacuum-energy density as a negative quantity. In the equation of state [Equation (5.28)], ρ_λ and p_λ are then scale-independent constants.

Early Time Dependence. It follows from the above scale dependences that the curvature term in Equation (5.17) obeys the following inequality in the limit of small a :

$$\frac{kc^2}{a^2} \ll \frac{8\pi G}{3} \rho + \frac{\lambda}{3}. \quad (5.35)$$

In fact, this inequality is always true when

$$k = +1, \quad p > -\frac{1}{3}\rho c^2, \quad w > -\frac{1}{3}, \quad \lambda > 0. \quad (5.36)$$

Then we can neglect the curvature term and the λ term in Equation (5.17), which simplifies to

$$\frac{\dot{a}}{a} = H(t) = \left(\frac{8\pi G}{3} \rho \right)^{1/2} \propto a^{-3(1+w)/2}. \quad (5.37)$$

Let us now find the time dependence of a by integrating this differential equation:

$$\int da a^{-1+3(1+w)/2} \propto \int dt,$$

to obtain the solutions

$$a^{3(1+w)/2} \propto t \quad \text{for } w \neq -1, \quad \ln a \propto t \quad \text{for } w = -1.$$

Solving for a ,

$$a(t) \propto t^{2/3(1+w)} \quad \text{for } w \neq -1, \quad a(t) \propto e^{\text{const.} \cdot t} \quad \text{for } w = -1. \quad (5.38)$$

In the two epochs of matter domination and radiation domination we know the value of w . Inserting this we obtain the time dependence of a for a matter-dominated universe,

$$a(t) \propto t^{2/3}, \quad (5.39)$$

and for a radiation-dominated universe,

$$a(t) \propto t^{1/2}. \quad (5.40)$$

Big Bang. We find the starting value of the scale of the Universe independently of the value of k in the curvature term neglected above:

$$\lim_{t \rightarrow 0} a(t) = 0. \quad (5.41)$$

In the same limit the rate of change \dot{a} is obtained from Equation (5.37) with any w obeying $w > -1$:

$$\lim_{t \rightarrow 0} \dot{a}(t) = \lim_{t \rightarrow 0} a^{-1}(t) = \infty. \quad (5.42)$$

It follows from Equations (5.32) and (5.33) that an early radiation-dominated Universe was characterized by extreme density and pressure:

$$\begin{aligned} \lim_{t \rightarrow 0} \rho_r(t) &= \lim_{t \rightarrow 0} a^{-4}(t) = \infty, \\ \lim_{t \rightarrow 0} p_r(t) &= \lim_{t \rightarrow 0} a^{-4}(t) = \infty. \end{aligned}$$

In fact, these limits also hold for any w obeying $w > -1$.

Actually, we do not even need an equation of state to arrive at these limits. Provided $\rho c^2 + 3p$ was always positive and λ negligible, we can see from Equations (5.6)

and (5.18) that the Universe has always decelerated. It then follows that a must have been zero at some time in the past. Whether Friedmann's equations can in fact be trusted to that limit is another story which we shall come back to later. The time $t = 0$ was sarcastically called the *Big Bang* by Fred Hoyle, who did not like the idea of an expanding Universe starting from a singularity, but the name has stuck. Since about 1988 the steady state theory has been abandoned because of the discovery of early quasars.

Late Einstein–de Sitter Evolution. The conclusions we derived from Equation (5.35) were true for past times in the limit of small a . However, the recent evolution and the future depend on the value of k and on the value of λ . For $k = 0$ and $k = -1$ the expansion always continues, following Equation (5.38), and a positive value of λ boosts the expansion further.

In a matter-dominated Einstein–de Sitter universe which is flat and has $\Omega_\lambda = 0$, Friedmann's Equation (5.4) can be integrated to give

$$t(z) = \frac{2}{3H_0}(1+z)^{-3/2}, \quad (5.43)$$

and the present age of the Universe at $z = 0$ would be

$$t_0 = \frac{2}{3H_0}. \quad (5.44)$$

In that case the size of the Universe would be $ct_0 = 2h^{-1}$ Gpc. Inserting the value of H_0 used in Equation (1.21), $H_0 = 0.696 \text{ km s}^{-1} \text{ Mpc}^{-1}$, one finds

$$t_0 = 9.27 \text{ Gyr}. \quad (5.45)$$

This is in obvious conflict with t_0 as determined from the ages of the oldest known star in the Galaxy in Equation (1.24), 13.5 ± 2.9 Gyr. Thus the flat-universe model with $\Omega_\lambda = 0$ is in trouble.

Evolution of a Closed Universe. In a closed matter-dominated universe with $k = +1$ and $\lambda = 0$, the curvature term kc^2/a^2 drops with the second power of a , while, according to Equation (5.30), the density drops with the third power, so the inequality [Equation (5.35)] is finally violated. This happens at a scale a_{max} such that

$$a_{\text{max}}^{-2} = \frac{8\pi G\rho_m}{3c^2}, \quad (5.46)$$

and the expansion halts because $\dot{a} = 0$ in Equation (5.4). Let us call this the *turnover time* t_{max} . At later times the expansion turns into contraction, and the Universe returns to zero size at time $2t_{\text{max}}$. That time is usually called the *Big Crunch*. For $k = +1$ Friedmann's Equation (5.4) then takes the form

$$\frac{da}{dt} = \sqrt{\frac{8\pi}{3}G\rho_m(a)a^2 - c^2}.$$

Then t_{\max} is obtained by integrating t from 0 to t_{\max} and a from 0 to a_{\max} ,

$$t_{\max} = \frac{1}{c} \int_0^{a_{\max}} da \left(\frac{8\pi G}{3c^2} \rho_m(a) a^2 - 1 \right)^{-1/2}. \quad (5.47)$$

To solve the a integral we need to know the energy density $\rho_m(a)$ in terms of the scale factor, and we need to know a_{\max} . Let us take the mass of the Universe to be M . We have already found in Equation (2.43) that the volume of a closed universe with Robertson–Walker metric is

$$V = 2\pi^2 a^3.$$

Since the energy density in a matter-dominated universe is mostly pressureless dust,

$$\rho_m = \frac{M}{V} = \frac{M}{2\pi^2 a^3}. \quad (5.48)$$

This agrees perfectly with the result [Equation (5.30)] that the density is inversely proportional to a^3 . Obviously, the missing proportionality factor in Equation (5.30) is then $M/2\pi^2$. Inserting the density [Equation (5.48)] with $a = a_{\max}$ into Equation (5.46) we obtain

$$a_{\max} = \frac{4MG}{3\pi c^2}. \quad (5.49)$$

We can now complete the integral in Equation (5.47):

$$t_{\max} = \frac{\pi}{2c} a_{\max} = \frac{2MG}{3c^3}. \quad (5.50)$$

Although we might not know whether we live in a closed universe, we certainly know from the ongoing expansion that $t_{\max} > t_0$. Using the value for t_0 from Equation (1.21) we find a lower limit to the mass of the Universe:

$$M > \frac{3t_0 c^3}{2G} = 1.30 \times 10^{23} M_{\odot}. \quad (5.51)$$

Actually, the total mass inside the present horizon is estimated to be about $10^{22} M_{\odot}$.

The dependence of t_{\max} on Ω_m can also be obtained:

$$t_{\max} = \frac{\pi \Omega_m}{2H_0(\Omega_m - 1)^{3/2}}. \quad (5.52)$$

The three cases $k = -1, 0, +1$ with $\lambda = 0$ are illustrated qualitatively in Figure 5.1. All models have to be consistent with the scale and rate of expansion today, $a = 1$ and \dot{a}_0 , at time t_0 . Following the curves back in time one notices that they intersect the time axis at different times. Thus what may be called time $t = 0$ is more recent in a flat universe than in an open universe, and in a closed universe it is even more recent.

The Radius of the Universe. The spatial curvature is given by the Ricci scalar R introduced in Equation (3.18), and it can be expressed in terms of Ω :

$$R = 6H^2(\Omega - 1). \quad (5.53)$$

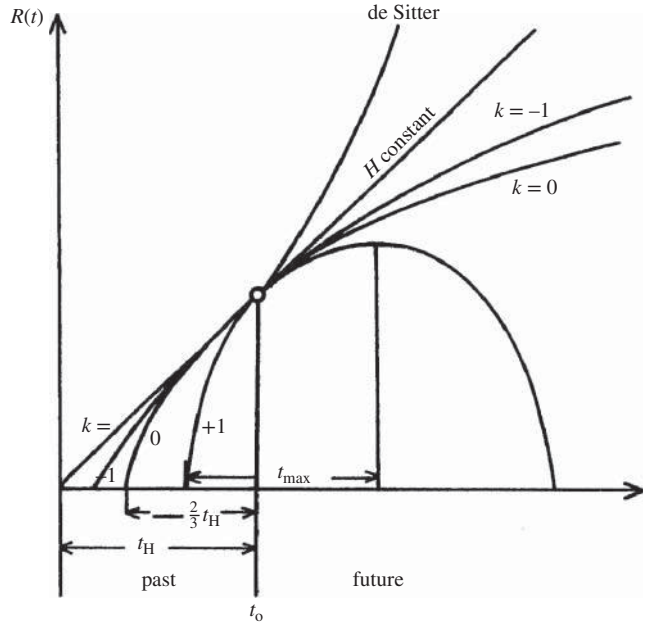


Figure 5.1 Time dependence of the cosmic scale $R(t) = a(t)$ in various scenarios, all of which correspond to the same constant slope $H = H_0$ at the present time t_0 . $k = +1$: a closed universe with a total lifetime $2t_{\max}$. It started more recently than a flat universe would have. $k = 0$: a flat universe which started $\frac{2}{3}t_H$ ago. $k = -1$: an open universe which started at a time $\frac{2}{3}t_H < t < t_H$ before the present time. de Sitter: an exponential (inflationary) scenario corresponding to a large cosmological constant. This is also called the Lemaitre cosmology.

Obviously, R vanishes in a flat universe, and it is only meaningful when it is non-negative, as in a closed universe. It is conventional to define a ‘radius of curvature’ that is also valid for open universes:

$$r_U \equiv \sqrt{\frac{6}{R}} = \frac{1}{H\sqrt{|\Omega - 1|}}. \quad (5.54)$$

For a closed universe, r_U has the physical meaning of the radius of a sphere.

Late Friedmann–Lemaître Evolution. When $\lambda > 0$, the recent past and the future take an entirely different course (we do not consider the case $\lambda < 0$, which is of mathematical interest only). Since ρ_λ and Ω_λ are then scale-independent constants, they will start to dominate over the matter term and the radiation term when the expansion has reached a given scale. Friedmann’s Equation (5.18) can then be written

$$\frac{2\ddot{a}}{a} = 3H_0^2\Omega_\lambda.$$

From this one sees that the expansion will accelerate regardless of the value of k . In particular, a closed universe with $k = +1$ will ultimately not contract, but expand at an accelerating pace.

Let us now return to the general expression [Equation (5.14)] for the normalized age $t(z)/t_0$ or $t(a)/t_0$ of a universe characterized by k and energy density components Ω_m , Ω_r and Ω_λ . Inserting the Ω components into Equations (5.13) and (5.14) we have

$$\frac{\dot{a}^2}{a^2} = H^2(t) = H_0^2[(1 - \Omega_0)a^{-2} + \Omega_m(a) + \Omega_r(a) + \Omega_\lambda(a)],$$

or

$$t(z) = \frac{1}{H_0} \int_0^{1/(1+z)} da [(1 - \Omega_0) + \Omega_m a^{-1} + \Omega_r a^{-2} + \Omega_\lambda a^2]^{-1/2}. \quad (5.55)$$

The integral can easily be carried out analytically when $\Omega_\lambda = 0$. But this is now of only academic interest, since we know today that $\Omega_\lambda \approx 0.7$ as we shall see later. Thus the integral is best solved numerically (or analytically in terms of hypergeometric functions or the Weierstrass modular functions [3, 4]).

The lookback time is given by the same integral with the lower integration limit at $1/(1+z)$ and the upper limit at 1. The proper distance [Equation (2.39)] is then

$$d_p(z) = \chi(z) = ct(z). \quad (5.56)$$

In Figure 5.2 we plot the lookback time $t(z)/t_0$ and the age of the Universe $1 - t(z)/t_0$ in units of t_0 as functions of redshift for the parameter values $\Omega_m = 0.27$, $\Omega_\lambda = 1 - \Omega_m$. At infinite redshift the lookback time is unity and the age of the Universe is zero.

Another important piece of information is that $\Omega_0 \approx 1.0$ (Table A.6). The vacuum term [Equation (5.8)] (almost) vanishes, in which case we can conclude that the geometry of our Universe is (almost) flat. With $\Omega_0 = 1.0$ and Ω_r well known, the integral [Equation (5.55)] really depends on only one unknown parameter, $\Omega_m = 1 - \Omega_\lambda$.

From the values $\Omega_\lambda \approx 0.7$ and $\Omega_m \approx 1 - 0.7 = 0.3$, one can conclude that the cosmological constant has already been dominating the expansion for some time. We shall come back to this later.

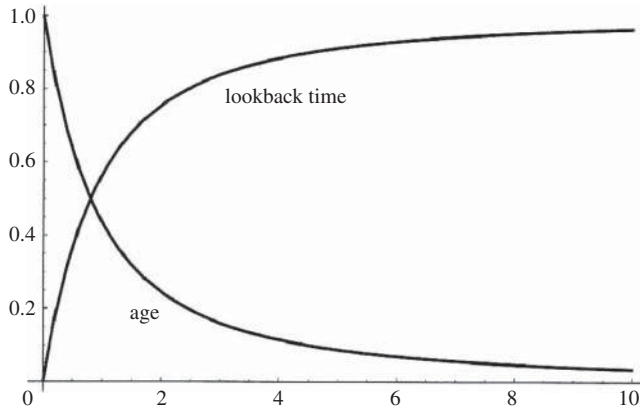


Figure 5.2 The lookback time and the age of the Universe normalized to t_0 as functions of redshift for the parameter values $\Omega_m = 0.27$, $\Omega_\lambda = 1 - \Omega_m$.

5.2 de Sitter Cosmology

Let us now turn to another special case for which the Einstein equation can be solved. Consider a homogeneous flat universe with the Robertson–Walker metric in which the density of pressureless dust is constant, $\rho(t) = \rho_0$. Friedmann’s Equation (5.17) for the rate of expansion including the cosmological constant then takes the form

$$\frac{\dot{a}(t)}{a(t)} = H, \quad (5.57)$$

where H is now a constant:

$$H = \sqrt{\frac{8\pi}{3}G\rho_0 + \frac{\lambda}{3}}. \quad (5.58)$$

This is clearly true when $k = 0$ but is even true for $k \neq 0$: since the density is constant and a increases without limit, the curvature term kc^2/R^2 will eventually be negligible. The solution to Equation (5.57) is obviously an exponentially expanding universe:

$$a(t) \propto e^{Ht}. \quad (5.59)$$

This is drawn as the de Sitter curve in Figure 5.1. Substituting this function into the Robertson–Walker metric [Equation (2.32)] we obtain the de Sitter metric

$$ds^2 = c^2 dt^2 - e^{2Ht}(dr^2 + r^2 d\theta^2 + r^2 \sin^2\theta d\phi^2) \quad (5.60)$$

with r replacing σ . In 1917 de Sitter published such a solution, setting $\rho = p = 0$, thus relating H directly to the cosmological constant λ . The same solution of course follows even with $\lambda = 0$ if the density of dust ρ is constant. Eddington characterized the *de Sitter universe* as ‘motion without matter’, in contrast to the static *Einstein universe* that was ‘matter without motion’.

If one introduces two test particles into this empty de Sitter universe, they will appear to recede from each other exponentially. The force driving the test particles apart is very strange. Let us suppose that they are at spatial distance ra from each other, and that λ is positive. Then the equation of relative motion of the test particles is given by Equation (5.5) including the λ term:

$$\frac{d^2(ra)}{dt^2} = \frac{\lambda}{3}ra - \frac{4\pi}{3}G(\rho + 3pc^{-2})ra. \quad (5.61)$$

The second term on the right-hand side is the decelerating force due to the ordinary gravitational interaction. The first term, however, is a force due to the vacuum-energy density, proportional to the distance r between the particles!

If λ is positive as in the Einstein universe, the force is repulsive, accelerating the expansion. If λ is negative, the force is attractive, decelerating the expansion just like ordinary gravitation. This is called an *anti-de Sitter* universe. Since λ is so small [see Equation (5.22)] this force will only be of importance to systems with mass densities of the order of the vacuum energy. The only known systems with such low densities are the large-scale structures, or the full horizon volume of cosmic size. This is the reason for the name *cosmological constant*. In a later chapter we shall meet inflationary universes with exponential expansion.

Although the world is not devoid of matter and the cosmological constant is small, the de Sitter universe may still be of more than academic interest in situations when ρ changes much more slowly than the scale a . The de Sitter metric then takes the form

$$ds^2 = (1 - r^2 H^2) dt^2 - (1 - r^2 H^2)^{-1} dr^2 - r^2(d\theta^2 + \sin^2\theta d\phi^2). \quad (5.62)$$

Not that the coordinates here violate the cosmological principle. There is an *inside* region in the de Sitter space at $r < H^{-1}$, for which the metric tensor component g_{00} is positive and g_{11} is negative. This resembles the region *outside* a black hole of Schwarzschild radius $r_c = H^{-1}$, at $r > r_c$, where g_{00} is positive and g_{11} is negative. Outside the radius $r = H^{-1}$ in de Sitter space and inside the Schwarzschild black hole these components of the metric tensor change sign.

At $r = H^{-1}$, $g_{00} = 1 - r^2 H^2$ vanishes and $g_{11} = -(1 - r^2 H^2)^{-1}$ is singular. In the sub-space defined by $0 \leq r \leq H^{-1}$ we can transform the singularity away by the substitution $u^2 = H^{-1} - r$. The new radial coordinate, u , is then in the range $0 \leq u \leq H^{-1}$ and the nonsingular metric becomes

$$ds^2 = (1 - r^2 H^2) dt^2 - 4H^{-1}(1 + rH)^{-1} du^2 - r^2(d\theta^2 + \sin^2\theta d\phi^2). \quad (5.63)$$

From this example we see that some singularities may be just the consequence of a badly chosen metric and not a genuine property of the theory.

The interpretation of this geometry is that the de Sitter metric describes an expanding space-time surrounded by a black hole. Inside the region $r = H^{-1}$ no signal can be received from distances outside H^{-1} because there the metric corresponds to the inside of a black hole! In an anti-de Sitter universe the constant attraction ultimately dominates, so that the expansion turns into contraction. Thus de Sitter universes are open and anti-de Sitter universes are closed.

Let us study the particle horizon r_H in a de Sitter universe. Recall that this is defined as the location of the most distant visible object, and that the light from it started on its journey towards us at time t_H . From Equation (2.48) the particle horizon is at

$$r_H(t) = R(t)\chi_{\text{ph}} = R(t) \int_{t_H}^{t_0} \frac{dr'}{R(t')}. \quad (5.64)$$

Let us choose t_H as the origin of time, $t_H = 0$. The distance $r_H(t)$ as a function of the time of observation t then becomes

$$r_H(t) = H^{-1} e^{Ht} (1 - e^{-Ht}). \quad (5.65)$$

The comoving distance to the particle horizon, χ_{ph} , quickly approaches the constant value H^{-1} . Thus for a comoving observer in this world the particle horizon would always be located at H^{-1} . Points which were inside this horizon at some time will be able to exchange signals, but events outside the horizon cannot influence anything inside this world.

The situation in a Friedmann universe is quite different. There the time dependence of a is not an exponential but a power of t , Equation (5.38), so that the comoving distance χ_{ph} is an increasing function of the time of observation, not a constant. Thus points which were once at space-like distances, prohibited to exchange signals with each other, will be causally connected later, as one sees in Figure 2.1.

The importance of the de Sitter model will be illustrated later when we deal with exponential expansion at very early times in inflationary scenarios.

5.3 The Schwarzschild Model

If the Einstein Equations (3.29) were difficult to derive, it was even more difficult to find solutions to this system of ten coupled nonlinear differential equations. A particularly simple case, however, is a single spherically symmetric star of mass M surrounded by a vacuum universe, or in any case far away from the gravitational influence of other bodies.

The Schwarzschild Metric. Let the metric near the star be described by a time coordinate t and a radial elevation r measured from the center of the star. We assume stable conditions so that the field is stationary, but it varies with elevation. The metric is then not flat, but the 00 time–time component and the 11 space–space component must be modified by some functions of r . Thus it can be written in the form

$$ds^2 = B(r)c^2 dt^2 - A(r) dr^2, \tag{5.66}$$

where $B(r)$ and $A(r)$ have to be found by solving the Einstein equations.

Far away from the star the space-time may be taken to be flat. This gives us the asymptotic conditions

$$\lim_{r \rightarrow \infty} A(r) = \lim_{r \rightarrow \infty} B(r) = 1. \tag{5.67}$$

From Equation (3.41) the Newtonian limit of g_{00} is known. Here $B(r)$ plays the role of g_{00} ; thus we have

$$B(r) = 1 - \frac{2GM}{c^2 r}. \tag{5.68}$$

To obtain $A(r)$ from the Einstein equations is more difficult, and we shall not go to the trouble of deriving it. The exact solution found by *Karl Schwarzschild* (1873–1916) in 1916 preceded any solution found by Einstein himself. The result is simply

$$A(r) = B(r)^{-1}. \tag{5.69}$$

These functions clearly satisfy the asymptotic conditions [Equation (5.67)].

Let us introduce the concept of *Schwarzschild radius* r_c for a star of mass M , defined by $B(r_c) = 0$. It follows that

$$r_c \equiv \frac{2GM}{c^2}. \tag{5.70}$$

The physical meaning of r_c is the following. Consider a test body of mass m and radial velocity v attempting to escape from the gravitational field of the star. To succeed, its kinetic energy must overcome the gravitational potential. In the nonrelativistic case the condition for this is

$$\frac{1}{2}mv^2 \geq GMm/r. \tag{5.71}$$

The larger the ratio M/r of the star, the higher is the required escape velocity. Ultimately, in the ultra-relativistic case when $v = c$, only light can escape. At that point a nonrelativistic treatment is no longer justified. Nevertheless, it just so happens that the equality in Equation (5.71) fixes the radius of the star correctly to be precisely r_c , as defined above. Because nothing can escape the interior of r_c , not even light, *John A. Wheeler* coined the term black hole for it in 1967. Note that the escape velocity of objects on Earth is 11 km s^{-1} , on the Sun it is $2.2 \times 10^6 \text{ km h}^{-1}$, but on a black hole it is c .

This is the simplest kind of a *Schwarzschild black hole*, and r_c defines its event horizon. Inserting r_c into the functions A and B , the *Schwarzschild metric* becomes

$$d\tau^2 = \left(1 - \frac{r_c}{r}\right) dt^2 - \left(1 - \frac{r_c}{r}\right)^{-1} \frac{dr^2}{c^2}. \quad (5.72)$$

Note that The Schwarzschild metric resembles the de Sitter metric [Equation (5.62)]. In the Schwarzschild metric the coefficient of dt^2 is singular at $r = 0$, whereas the coefficient of dr^2 is singular at $r = r_c$. However, if we make the transformation from the radial coordinate r to a new coordinate u defined by

$$u^2 = r - r_c,$$

the Schwarzschild metric becomes

$$d\tau^2 = \frac{u^2}{u^2 + r_c} dt^2 - 4(u^2 + r_c) du^2.$$

The coefficient of dt^2 is still singular at $u^2 = -r_c$, which corresponds to $r = 0$, but the coefficient of du^2 is now regular at $u^2 = 0$.

5.4 Black Holes

A particularly fascinating and important case is a black hole, a star of extremely high density. Black holes are certainly the most spectacular prediction of general relativity, and they appear to be ubiquitous in the nuclei of bright and active galaxies.

The Schwarzschild metric has very fascinating consequences. Consider a spacecraft approaching a black hole with apparent velocity $v = dr/dt$ in the fixed frame of an outside observer. Light signals from the spacecraft travel on the light cone, $d\tau = 0$, so that

$$\frac{dr}{dt} = c \left(1 - \frac{r_c}{r}\right). \quad (5.73)$$

Thus the spacecraft appears to slow down with decreasing r , finally coming to a full stop as it reaches the mathematical singularity of dt at the *event horizon* $r = r_c$ in the expression

$$c dt = \frac{dr}{1 - r_c/r}. \quad (5.74)$$

The time intervals dt between successive crests in the wave of the emitted light become longer, reaching infinite wavelength at the singularity. Although the velocity

of the emitted photons is unchanged c their frequency ν goes to zero, and the energy $E = h\nu$ of the signal vanishes. One cannot receive signals from beyond the event horizon because photons cannot have negative energy. Thus the outside observer sees the spacecraft slowing down and the signals redshifting until they cease completely.

The pilot in the spacecraft using local coordinates sees the passage into the black hole entirely differently. If he started out at distance r_0 with velocity $dr/dt = 0$ at time t_0 , he will have reached position r at proper time τ , which we can find by integrating $d\tau$ in Equation (5.72) from 0 to τ :

$$\int_0^\tau \sqrt{d\tau^2} = \tau = \int_{r_0}^r \left[\frac{1 - r_c/r}{(dr/dt)^2} - \frac{1}{c^2(1 - r_c/r)} \right]^{1/2} dr. \quad (5.75)$$

The result depends on $dr(t)/dt$, which can only be obtained from the equation of motion. The pilot considers that he can use Newtonian mechanics, so he may take

$$\frac{dr}{dt} = c \sqrt{\frac{r_c}{r}}.$$

The result is then (Problem 11.4):

$$\tau \propto (r_0 - r)^{3/2}. \quad (5.76)$$

However, many other expressions for $dr(t)/dt$ also make the integral in Equation (5.75) converge.

Thus the singularity at r_c does not exist to the pilot, his comoving clock shows finite time when he reaches the event horizon. The fact that the singularity at r_c does not exist in the local frame of the spaceship indicates that the horizon at r_c is a mathematical singularity and not a physical singularity. The singularity at the horizon arises because we are using, in a region of extreme curvature, coordinates most appropriate for flat or mildly curved space-time. Alternate coordinates, more appropriate for the region of a black hole and in which the horizon does not appear as a singularity, were invented by Eddington (1924) and rediscovered by Finkelstein (1958). If we were able to observe the collapse of a neutron star towards the Schwarzschild radius into a black hole it would appear to take a very long time. Towards the end of it, the ever-redshifting light would fade and finally disappear completely.

Note from the metric Equation (5.72) that inside r_c the time term becomes negative and the space term positive, thus space becomes timelike and time spacelike. The implications of this are best understood if one considers the shape of the light cone of the spacecraft during its voyage (see Figure 5.3). Outside the event horizon the future light cone contains the outside observer, who receives signals from the spacecraft. Nearer r_c the light cone narrows and the slope dr/dt steepens because of the approaching singularity in expression on the the right-hand side of Equation (5.73). The portion of the future space-time which can receive signals therefore diminishes.

Since the time and space axes have exchanged positions inside the horizon, the future light cone is turned inwards and no part of the outside space-time is included in the future light cone. The slope of the light cone is vertical at the horizon. Thus it defines, at the same time, a cone of zero opening angle around the original time axis, and a cone of 180° around the final time axis, encompassing the full space-time

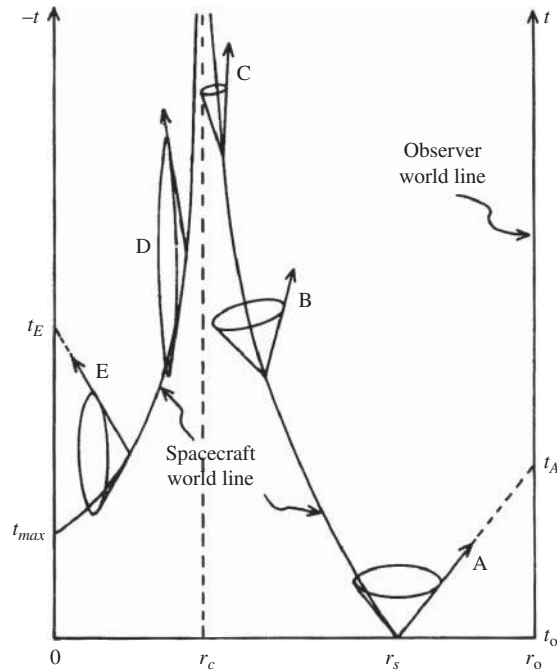


Figure 5.3 The world line of a spacecraft falling into a Schwarzschild black hole. A, the journey starts at time t_0 when the spacecraft is at a radius r_s , far outside the Schwarzschild radius r_c , and the observer is at r_o . A light signal from the spacecraft reaches the observer at time $t_A > t_0$ (read time on the right-hand vertical scale!). B, nearer the black hole the future light cone of the spacecraft tilts inward. A light signal along the arrow will still reach the observer at a time $t_B \gg t_A$. C, near the Schwarzschild radius the light cone narrows considerably, and a light signal along the arrow reaches the observer only in a very distant future. D, inside r_c the time and space directions are interchanged, time running from up to down on the left-hand vertical scale. All light signals will reach the center of the hole at $r = 0$, and none will reach the observer. The arrow points in the backward direction, so a light signal will reach the center after the spacecraft. E, the arrow points in the forward direction of the hole, so that a light signal will reach the center at time t_E , which is earlier than t_{max} , when the spacecraft ends its journey.

of the black hole. As the spacecraft approaches the center, dt/dr decreases, defining a narrowing opening angle which always contains the center.

When the center is reached the metric [Equation (5.72)] has a physical singularity. the spacecraft no longer has a future. One cannot define field equations there, so general relativity breaks down, unable to predict what will happen. Once across r_c the spacecraft reaches the center of the black hole rapidly. For a hole of mass $10M_\odot$ this final passage lasts about 10^{-4} s.

Some people have speculated that matter or radiation falling in might ‘tunnel’ through a ‘wormhole’ out into another universe. Needless to say, all such ideas are purely theoretical speculations with no hope of experimental verification.

An interesting quantity is the Schwarzschild radius of the Universe, $r_{c,U}$. Combining Equations (5.50) and (5.70) we find

$$r_{c,U} = 3ct_{\max} > 12 \text{ Gpc.} \quad (5.77)$$

Comparing this number with the much smaller Hubble radius $3h^{-1}$ Gpc in Equation (1.14) we might conclude that we live inside a black hole! However, the Schwarzschild metric is static and takes the black hole to be surrounded by a nonexpanding universe, whereas the Hubble radius recedes in expanding Friedmann models with superluminal velocity, as was seen in Equation (2.52), so it will catch up with $r_{c,U}$ at some time. Actually, it makes more sense to describe the Big Bang singularity as a *white hole*, which is a time-reversed black hole. A white hole only emits and does not absorb. It has a horizon over which nothing gets in, but signals from inside do get out.

Black holes possessing either charge or angular momentum are called *Reissner–Nordström black holes* and *Kerr–Newmann black holes*, respectively, and they are described by different metrics. It is natural to consider that matter attracted by a hole has angular momentum. Matter can circulate a hole in stable orbits with radii exceeding $3r_c$, but if it comes any closer it starts to spiral in towards the horizon, and is soon lost into the hole with no possibility to escape. Since angular momentum is conserved, the infalling matter must speed up the rotation of the hole. However, centrifugal forces set a limit on the angular momentum J that a rotating black hole can possess:

$$J \leq \frac{GM^2}{c}. \quad (5.78)$$

This does not imply that the hole is ripped into pieces with one increment of rotating matter, rather, that it could never have formed in the first place. Remember that angular momentum is energy, and energy is curvature, so incremental energy is modifying the space-time geometry of the black hole, leading to a smaller event horizon. Thus the angular momentum can never overcompensate the gravitational binding energy. If it could, there would be no event horizon and we would have the case of a visible singularity, also called a *naked singularity*. Since nobody has conceived of what such an object would look like, *Stephen Hawking* and *Roger Penrose* have conjectured that space-time singularities should always be shielded from inspection by an event horizon. This is called the principle of *cosmic censorship*—in Penrose’s words ‘Nature abhors a naked singularity’. The reader might find further enjoyment reading the book by Hawking and Penrose on this subject [5].

Event Horizons. Classically a black hole is a region from which nothing can escape, not even light. Black holes are mathematically simple objects obeying general relativity, as seen from outside their event horizon, they have only the three properties: mass, electric charge and angular momentum. Their size depends only on their mass so that all holes with the same mass are identical and exactly spherical, unless they rotate. All other properties possessed by stars, such as shape, solid surface, electric dipole moment, magnetic moments, as well as any detailed outward structure, are absent. This has led to John Wheeler’s famous statement ‘black holes have no hair’.

However, there may be reasons to doubt the properties of classical event horizons. The Schwarzschild metric in Equation (5.72) is perfectly classical, the radial coordinate and the time are exact. But this ignores the fact that coordinates in quantum mechanics are uncertain, merely statistical quantities subject to fluctuations. If an object meets the event horizon exactly at radius r , the time for this occurrence is completely undetermined according to Heisenberg's uncertainty relation. Thus the left hand of the infalling pilot may disappear from the view of the outside observer a week before his right hand which is still seen waving goodbye.

Quantum theory appears to dictate that the event horizon actually is a highly energetic region, or 'firewall', that would burn the astronaut to a crisp, in contradiction with the property endowed black holes by general relativity, that they 'have no hair'. To really understand this *black-hole firewall paradox* would require a theory encompassing quantum mechanics and general relativity, still missing.

Instead of a firewall, Hawking has recently speculated in a brief paper without quantitative calculations [6], that quantum mechanics and general relativity remain intact, but black holes are not surrounded by a classical event horizon or a firewall hindering radiation and information to escape to infinity. Rather they are surrounded by a much more benign 'apparent horizon' which would enable the release of some energy and information temporarily.

Such apparent horizons would persist only for a period of time. This suggests that black holes are collapsed objects which should be redefined as metastable bound states. Thus black holes would not be black except for a period of time. Inside the event horizon, the metric and matter fields would be classical and deterministic but chaotic. The key to his claim is that quantum effects around the black hole cause space-time to fluctuate too wildly for a sharp boundary surface to exist [6].

In general relativity, for an unchanging black hole, these two horizons are identical, because light trying to escape from inside a black hole can reach only as far as the event horizon and will be held there. However, the two horizons can, in principle, be distinguished. If more matter gets swallowed by the black hole, its event horizon will swell and grow larger than the apparent horizon. Conversely, in the 1970s, Hawking also showed that black holes can slowly shrink, spewing out *Hawking radiation* (see below). In that case, the event horizon would, in theory, become smaller than the apparent horizon. Hawking's new suggestion is that the apparent horizon is the real boundary. Unlike the event horizon, the apparent horizon can eventually dissolve. Anything in principle could then get out of a black hole. The new idea that there are no inside points from which you cannot escape a black hole is in some ways an even more radical and problematic suggestion than the existence of firewalls [6].

If Hawking is correct, there could even be no singularity at the core of the black hole. Instead, matter would be only temporarily held behind the apparent horizon, which would gradually move inward owing to the pull of the black hole, but would never quite crunch down to the center. Information about this matter would not be totally destroyed, but would be released in a vastly different form, making it almost impossible to work out what the swallowed objects once were like [6].

The *information loss paradox* can be stated neatly as follows: if one throws a book into a black hole, all that comes out is blackbody radiation or chaotic radiation, the

information contained in the book is lost forever. This is a paradox because quantum theory states that the information of the initial state should never disappear. In order to resolve this paradox it is necessary to construct microscopic states of the black hole and to give a statistical-mechanical explanation for the black hole entropy, which is difficult within general relativity because of the no-hair theorem. So far the paradox still remains since a complete description of an evaporating black hole has not yet been established.

In 1973 *J. Bekenstein* noted [7] that there are certain similarities between the size of the event horizon of a black hole and entropy. When a star has collapsed to the size of its Schwarzschild radius, its event horizon will never change (to an outside observer) although the collapse continues. Thus entropy s could be defined as the surface area A of the event horizon times some proportionality factor,

$$s = \frac{Akc^3}{4G\hbar}, \quad (5.79)$$

the *Bekenstein–Hawking formula*. For a spherically symmetric black hole of mass M the surface area is given by

$$A = 16\pi M^2 G^2 / c^4. \quad (5.80)$$

A can increase only if the black hole devours more mass from the outside, but A can never decrease because no mass will leave the horizon. In Hawking's parlance [6], however, this is true for the classical event horizon but not for the apparent horizon.

Inserting A into Equation (5.79), entropy comes out (in units of erg/K) proportional to M^2 :

$$s = M^2 4\pi kG / c\hbar. \quad (5.81)$$

Thus two black holes coalesced into one possess more entropy than they both had individually. This is illustrated in Figure 5.4.

Hawking Radiation. Stephen Hawking has shown [8, 9] that although no light can escape from black holes, they can nevertheless radiate if one takes quantum mechanics into account. It is a property of the *vacuum* that particle–antiparticle pairs such as e^-e^+ are continuously created out of nothing, to disappear in the next moment by *annihilation*, which is the inverse process. Since energy cannot be created or destroyed, one of the particles must have positive energy and the other one an equal amount of negative energy. They form a *virtual pair*, neither one is real in the sense that it could escape to infinity or be observed by us.

In a strong electromagnetic field the electron e^- and the positron e^+ may become separated by a Compton wavelength λ of the order of the Schwarzschild radius r_c . Hawking has previously shown that there is a small but finite probability for one of them to ‘tunnel’ through the barrier of the quantum vacuum and escape the black hole horizon as a real particle with positive energy, leaving the negative-energy particle inside the horizon of the hole. Since energy must be conserved, the hole loses mass in this process, a phenomenon called *Hawking radiation*. But, as noted before,

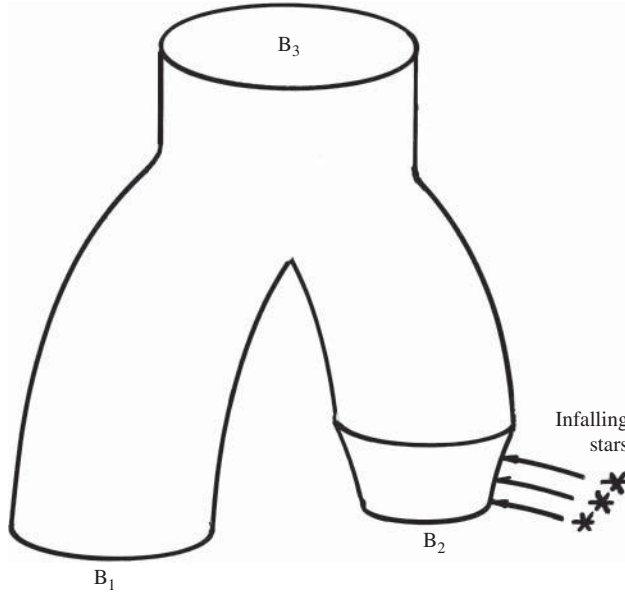


Figure 5.4 The merging of two black holes with event horizons B_1 and B_2 , respectively. The black hole 2 grows initially by swallowing infalling stars. The final black hole has an event horizon with event horizon $B_3 \geq B_1 + B_2$.

in quantum theory the radius r and the time t are conjugate quantities which must obey Heisenberg's uncertainty relation, so the conditions for Hawking radiation are uncertain.

The timescale of complete evaporation is

$$t \approx 10 \text{ Gyr} \left(\frac{M}{10^{12} \text{ kg}} \right)^3. \quad (5.82)$$

Thus small black holes evaporate fast, whereas heavy ones may have lifetimes exceeding the age of the Universe. Paradoxically, in some black holes the infall of matter may cause the black hole to evaporate faster than the horizon has time to form, as seen by a distant observer. This is a consequence of the fact that during gravitational collapse time dilation can increase without bound, sufficiently retarding implosion so as to prevent the formation of an event horizon and the potential loss of matter and information from the observable universe.

The analogy with entropy can be used even further. A system in thermal equilibrium is characterized by a unique temperature T throughout. When Hawking applied quantum theory to black holes, he found that the radiation emitted from particle-antiparticle creation at the event horizon is exactly thermal. The rate of particle emission is as if the hole were a black body with a unique temperature proportional to the gravitational field on the horizon, the *Hawking temperature*:

$$kT_H = \frac{c^3 h}{8\pi GM} = 6.15 \times 10^{-8} \frac{M_\odot}{M} \text{ K}. \quad (5.83)$$

Black Hole Creation. Primordial black holes may have been created in the early universe by initial homogeneities, inflation, phase transitions, bubble collisions, or the decay of cosmic loops. Small primordial black holes are not expected to exist because they may already have evaporated.

The neutrons in a neutron star form a cold Fermi gas in which the quantum degeneracy pressure of the neutrons prevent the star from collapse. If the mass of the star exceeds the *Tolman–Oppenheimer–Volkoff limit* of $2.0\text{--}3.0M_{\odot}$ it may not always be stabilized against bounce. They then follow the route of further gravitational collapse to become a hypothetical *quark star* or a black hole.

The fate of a collapsing spherical star can be illustrated schematically by a light cone diagram. Consider the evolution in time of an event horizon corresponding to the circular equatorial intersection of the star. With increasing time, vertically upwards in Figure 5.5, the the equatorial intersection shrinks and in consequence light rays

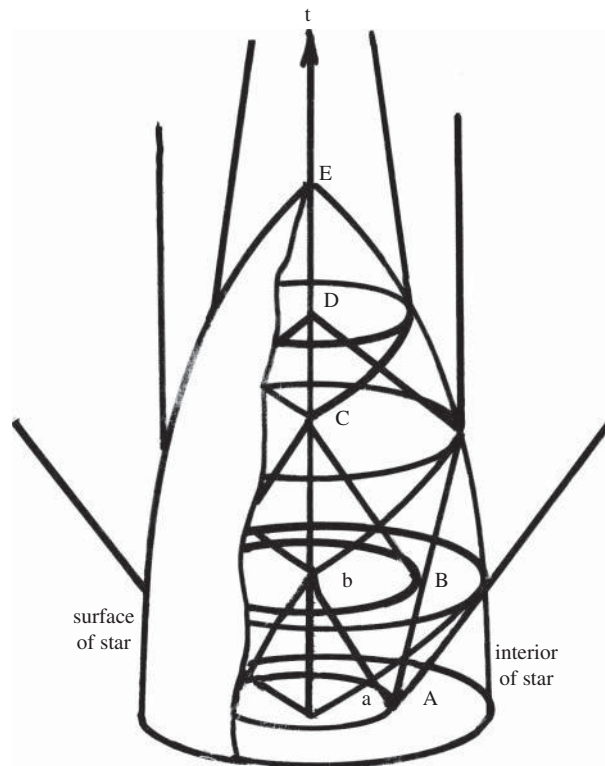


Figure 5.5 A space-time picture of the collapse of a spherical star to form a black hole. For an observer at the center of the surface (B), the surface (A) is his event horizon and (a) is his particle horizon. With increasing time, the radius of the star shrinks. For an observer at the center of the surface (C), the surface (B) is his event horizon and (b) is his particle horizon. At at the center of the surface (D) the particle horizon coincides with the event horizon (C). The surface (D) has become a trapped surface. At (E) no future light cone exists, the star has arrived at the singularity of the black hole.

toward the future steepen. When the star has shrunk to the size of the Schwarzschild radius, the equatorial section has become a trapped surface, the future of all light rays from it are directed towards the world line of the center of the star. For an outside observer the event horizon then remains eternally of the same size, but an insider would find that the trapped surface of the star is doomed to continue its collapse toward the singularity at zero radius. The singularity has no future light cone and cannot thus ever be observed (Figure 5.5).

Black holes are probably created naturally in the aging of stars. The collapse of an isolated heavy star is, however, not the only route to the formation of a black hole, and probably not even the most likely one. Binary stars are quite common, and a binary system consisting of a neutron star and a red giant is a very likely route to a black hole. The neutron star spirals in (see Figure 5.5), accretes matter from the red giant at a very high rate, about $1M_{\odot}$ per year, photons are trapped in the flow, gravity and friction heat the material in the accretion disc until it emits X-rays. When the temperature rises above 1 MeV neutrinos can carry off the thermal energy.

For example, Cyg X-1 is a black hole—a main-sequence star binary with a hole mass of more than about $10M_{\odot}$, probably the most massive black hole in a binary observed in the Galaxy. The enormous gravitational pull of the hole tears material from its companion star. This material then orbits the hole in a Saturnus-like accretion disc before disappearing into the hole. Finally, the main-sequence star explodes and becomes a neutron star, which will ultimately merge into the hole to form a heavier black hole.

Neutron stars and stellar black holes have masses of the order of $10M_{\odot}$, supermassive black holes have 10^6 – $10^{10}M_{\odot}$. There could be as many as 10^8 stellar black holes in our galaxy. Black holes weighing thousands or million times the Chandrasekhar limit of $1.44M_{\odot}$ cannot have been produced in the collapse of a single star. How they have been produced is currently not well understood. The most likely route is by accretion and coalescence of intermediate mass black holes.

Massive black holes are seen at very high redshift. The discovery of 10^9M_{\odot} black holes at $z > 7.0$ is mysterious since it is unclear how such massive objects could have formed and grown so quickly at early times.

Observations of Black Holes. Since black holes cannot be seen directly, one has to infer their existence from indirect evidence. Within our galaxy one can measure the orbits of neighboring individual stars. In this way one has deduced that the central few parsecs of our galaxy evidently hosts a dense and luminous star cluster and a very compact radio source Sgr A* [10]. Its intrinsic size is at most 10 light minutes which makes it the most likely candidate of a possible central black hole. There is evidence for little motion of the Sgr A* itself from the size and motions of the central compact radio source.

There are precise determinations of the elliptical orbits and the velocity vectors of about 30 stars in the vicinity, all having one common focal point located at the Sgr A*, see Figure 5.6. All the acceleration vectors intersect at Sgr A*, and the velocity vectors do not decrease with decreasing distance to Sgr A*, indicating that the stars move in

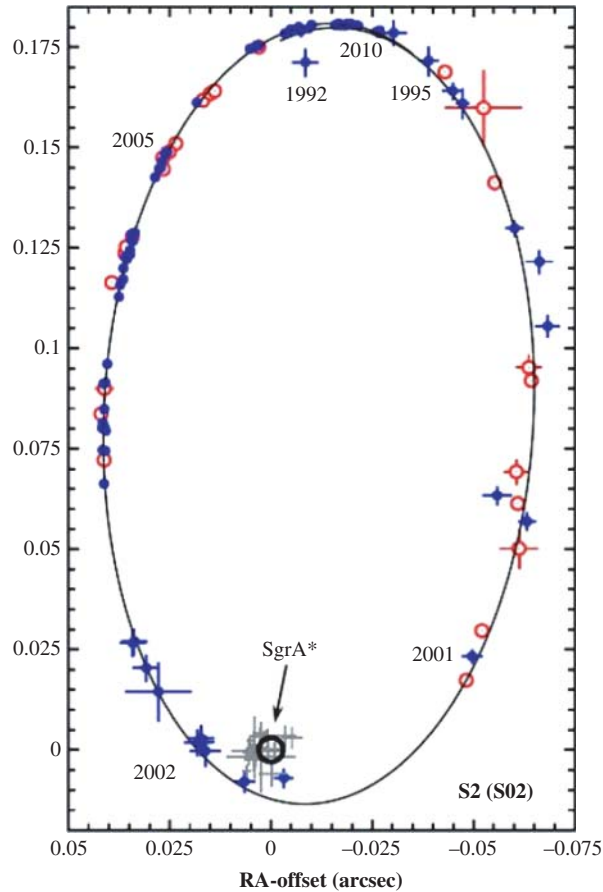


Figure 5.6 Orbit of the star S2 moving around Sgr A* [10]. Copyright 2010 by the American Physical Society. (See plate section for color version.)

a very strong gravitational field of a practically pointlike source of mass $4.4 \times 10^6 M_{\odot}$. The best measured orbit of the star S2 which has completed a full orbit since the beginning of 1992 is shown in Figure 5.6 [10].

There is a hot accreting zone around the event horizon where Sgr A* is accreting from the winds of surrounding stars at the outer boundary. Typically the radiated power from the accretion zone feeds the luminosity which is very weak both in radio and X-ray emission, but hot, thermal X-ray emission associated with Sgr A* has been claimed. It appears that there is a dynamically important magnetic field near the black hole, determining the structure of the accretion flow. If this field is accreted down to the event horizon it provides enough magnetic flux to explain the observed emission.

At the time of writing, an enigmatic small nonself-gravitating gas cloud is seen falling almost exactly onto Sgr A* at a nominal pericenter distance of only 2200 Schwarzschild radii ($r_c = 13^{\circ}$ km). Observations over the next decade should enlighten its fate and allow us to explore general relativity.

The accretion disks of black holes are typically truncated near the innermost circular orbit of matter, and hence the main effect is the deformation of a ring the size of this orbit bending near the black hole, rather than the effect of the event horizon itself. There will then be a shadow which is essentially a lensed image of the event horizon, and its shape is closely related to the closed orbit for photons around the black hole. Since photons that circle the black hole slightly within the photon orbit will end up inside the event horizon while photons just outside will escape to infinity, there is a rather sharp boundary between bright and dark regions. The shadow itself is thus indeed primarily a deficit of photons due to absorption by the event horizon.

A spectral turn-over indicates that the region becomes optically thin and allows us to see through to the event horizon. The shadow is in principle detectable with present-day technology and would allow many fundamental tests of general relativity and its alternatives.

Active galactic nuclei (AGN) furnish other indirect proof for black holes. AGN often exhibit quasars which require enormously powerful engines, much more than supernovae could furnish. It is now thought that every AGN is powered by a massive central black hole accreting radiation, matter and vacuum energy. The supermassive black hole Sgr A* is, however very inactive, so there is no quasar associated with it.

Spectacular flashes of energies 1000–10 000 larger than those of SNe (10^{54} erg, or the release of $1 M_{\odot}$ in a few 0.1 s) are seen homogeneously distributed in the universe. These *Gamma Ray Bursts* have a duration of 10^{-2} to 10^2 s. The very first one nearly started a third World War when United States intelligence suspected that the Soviet Union was testing a new arm above the Earth.

The GRBs are not ejected from black holes, rather they originate from the gravitational collapse of a neutron star to a black hole. One possible route for short GRBs is that they originate from an initial binary system consisting of a compact carbon–oxygen (CO) core star and a neutron star. The CO core explodes as a supernova, and part of its ejecta accretes onto the neutron star which reaches its critical mass and collapses to a black hole. A new neutron star is then generated by the supernova as a remnant. The energy released in the GRB is energetically dominant to supernovae, so it cannot originate in an SN. Long-duration GRBs are believed to be due to a continued energy injection mechanism which powers the forward shock, giving rise to an unusual and long-lasting afterglow.

If a neutron star and a black hole are gravitationally bound they can give rise to a new process of merging, resulting in a gravitational collapse leading to a new black hole with the emission of a GRB and possibly gravitational waves.

5.5 Extended Gravity Models

In Chapter 3.3 we derived the Einstein Equation (3.28) from the Einstein–Hilbert action

$$S = \int d^4x \sqrt{-g} \left[\frac{1}{2\kappa} R + \mathcal{L}_M(g_{\mu\nu}, \Psi) \right], \quad (5.84)$$

where Ψ represents some matter fields. Many proposals have been made to generalize this action. Suppose that the universe contains a scalar field ϕ not present in the Einstein equation. There are then different options how to introduce it, even the metric frame may need to be reformulated. Since the Einstein equation describes how matter affects the geometry and the curvature of space-time, the scalar field in the energy-momentum tensor will also affect the geometry.

The Robertson–Walker metric in Equation (2.32) determines the geometry of space-time by the components $g_{\mu\nu}$ in Equation (2.33), the affine connections $\Gamma_{\sigma\nu}^{\mu}$ in Equation (3.13), and the Ricci scalar R in (3.18) in what is called the *Einstein frame*. If the matter field ϕ is introduced into the action in such a way that no terms ϕR nor $F(\phi)R$ mixing matter and geometry appear, we are in the Einstein frame. The opposite case is the *Jordan frame*. The two frames are only mathematically different, because one can transform expressions in one frame into one in the other by a *conformal transformation* such as $\hat{g}_{\mu\nu} = \zeta^2(x)g_{\mu\nu}$, where $\zeta^2(x)$ is the conformal factor. There is no physical difference, only a choice of mathematical convenience.

Let us now rewrite the Einstein–Hilbert action in the most general form in the Jordan frame as

$$S_{total} = \int d^n x \sqrt{-g} \left[\frac{1}{2\kappa} f(R, \phi) + \mathcal{L}_{phi}(g_{\mu\nu}, \phi, \partial\phi) + \mathcal{L}_M(g_{\mu\nu}, \Psi) \right]. \quad (5.85)$$

Here we have replaced the curvature scalar R by a general function $f(R, \phi)$, and enlarged the dimensionality of space-time from 4 to n . The Lagrangian density is then usually of the form

$$\mathcal{L}_{phi} = -\frac{M^2}{2} \omega(\phi)(\partial\phi)^2 - V(\phi), \quad (5.86)$$

where we have denoted $(\partial\phi)^2 \equiv (\nabla^\alpha \phi)(\nabla_\alpha \phi)$. For a canonical scalar field the function ω equals unity.

Two simple versions of extended gravity are those which lack a scalar field ϕ but are nonlinear in R , and *scalar-tensor* theories. Some nonlinear functions $f(R)$ studied are of the form $(R + \alpha R^2)$ and $(R - \alpha/R)$. The dynamical equations are then written in the Einstein frame.

Scalar-tensor theories are of the form $f(R, \phi) = F(\phi)R$ in the Jordan frame, of which the classical (1961) *Brans–Dicke* theory with $F(\phi) = \phi$ is the simplest. Some other examples of scalar-tensor theories are $F(\phi) = \exp(-\phi)$ and $F(\phi) = \alpha\phi^2$.

If the theory has no scalar field ϕ one proceeds as for the Einstein–Hilbert action, varying the action S with respect to the metric $g_{\mu\nu}$ and setting the variation equal to zero. For scalar-tensor theories the equations of motion follow from setting the variation of S with respect to ϕ equal to zero, and next differentiating the component of the energy–momentum tensor $T_{\mu\nu}$ which depends on ϕ . After that, one again proceeds to vary the action S with respect to the metric $g_{\mu\nu}$ and setting the variation equal to zero.

Note that since the Ricci scalar R contains second derivatives of the metric, one ends up with fourth order terms of the metric which are difficult to analyze. A more tractable form of extended gravity can be obtained by using the *Palatini variation*. In this formulation the affine connection $\Gamma_{\sigma\nu}^{\mu}$ in Equation (3.13) is treated as an independent variable from the metric $g_{\mu\nu}$. Varying the action separately with respect to $g_{\mu\nu}$

and $\Gamma_{\sigma\nu}^{\mu}$ gives a system of second-order differential equations. This will of course give rise to somewhat different physics than the Einstein–Hilbert variation, so in due time we will probably have data permitting us to decide between them.

Another degree of freedom in the action [Equation (5.85)] is the dimensionality n of space-time. At this point in my treatise there is no obvious reason to assume that n is larger than four. However, the gravitational interaction is so much weaker than all the standard particle interactions, that one could imagine that it propagates in a different space-time. We shall come back to models in later chapters which explore this freedom.

Problems

1. On the solar surface the acceleration caused by the repulsion of a non-vanishing cosmological constant λ must be much inferior to the Newtonian attraction. Derive a limiting value of λ from this condition.
2. In Newtonian mechanics, the cosmological constant λ can be incorporated by adding to gravity an outward radial force on a body of mass m , a distance r from the origin, of $F = +m\lambda r/6$. Assuming that $\lambda = -10^{-20} \text{ yr}^{-2}$, and that F is the only force acting, estimate the maximum speed a body will attain if its orbit is comparable in size with the Solar System (0.5 light day [11]).
3. de Sitter's static universe has $a \propto e^{Ht}$, zero curvature of its co-moving coordinates ($k = 0$), and a proper density of all objects that is constant in time. Show that the co moving volume out to red-shift z is $V(z) = \frac{4}{3}\pi(cz/H)^3$, and hence that the number-count slope for objects at typical red-shift z becomes $[(3 + \alpha) \ln z]^{-1}$ for $z \gg 1$, where α is the spectral index for the objects [12].
4. Starting from Equation (5.55) with the parameters $\Omega_0 = 1$, $\Omega_r = 0$, show that the age of the Universe can be written

$$t_0 = \frac{2}{3H_0} \frac{\tanh^{-1} \sqrt{\Omega_\lambda}}{\sqrt{\Omega_\lambda}}.$$

5. A galaxy at $z = 0.9$ contains a quasar showing red-shift $z = 1.0$. Supposing that this additional red-shift of the quasar is caused by its proximity to a black hole, how many Schwarzschild radii away from the black hole does the light emitted by the quasar originate?
6. Estimate the gravitational red-shift z of light escaping from a galaxy of mass $10^9 M_\odot$ after being emitted from a nearby star at a radial distance of 1 kpc from the center of the galaxy. (Assume that all matter in the galaxy is contained within that distance [11].)
7. Light is emitted horizontally *in vacuo* near the Earth's surface, and falls freely under the action of gravity. Through what vertical distances has it fallen after travelling 1 km? Calculate the radial coordinate (expressed in Schwarzschild radii) at which light travels in a circular path around a body of mass M [11].

References

- [1] Rich, J. 2001 *Fundamentals of cosmology*. Springer.
- [2] Solà, J. 2001 *Nucl. Phys. B*, **95**, 29.
- [3] Cappi, A. 2001 *Astrophys. Lett. Commun.* **40**, 161.
- [4] Kraniotis, G. V. and Whitehouse, S. B. 2002 *Classical Quantum Gravity* **19**, 5073.
- [5] Hawking, S. and Penrose, R. 1996 *The nature of space and time*. Princeton University Press, Princeton, NJ.
- [6] Hawking, S. W. Preprint at <http://arxiv.org/abs/1401.5761> (2014).
- [7] Bekenstein, J. 1973 *Phys. Rev. D*, **7**, 2333.
- [8] Hawking, S. W. 1974 *Nature* **248**, 30.
- [9] Hawking, S. W. 1975 *Commun. Math. Phys.* **43**, 199.
- [10] Genzel, R., Eisenhauer, F. and Gillessen, S., 2010 *Rev. Mod. Phys* **82**, 3121.
- [11] Berry, M. V. 1989 *Principles of cosmology and gravitation*. Adam Hilger, Bristol.
- [12] Peacock, J. A. 1999 *Cosmological physics*. Cambridge University Press, Cambridge.

6

Thermal History of the Universe

This history should start at the Big Bang which we defined in the previous chapter as time zero. General relativity is a classical theory in which the coordinates t, x, y, z can take any real values. However, Equations (5.41) and (5.42) showed that at time zero no meaningful description of the Universe exists.

Actually there are problems already for small but nonzero times, at the energy scale where gravitational and quantum effects are of equal importance. Then we can not do particle physics neglecting gravitation nor can we describe the Universe with the Einstein equation neglecting quantum mechanics.

In Section 6.1 we begin with a definition of the Planck time, leaving the cosmic inflation for a later chapter. In Section 6.2 we continue with the primordial hot plasma which deals with particle physics symmetries, photons in thermal equilibrium, energy densities, relativistic versus nonrelativistic particles, spin and statistics, temperature and time scales.

We follow the cooling plasma until, in Section 6.3, photon and lepton decoupling occurs. Next comes photon reheating, neutrino decoupling and the recombination era. We end with a brief discussion on equilibrium theory.

In Section 6.4 we follow the thermal history of the nucleons for the momentous fusion processes in Big Bang nucleosynthesis (BBN) which has left us very important clues in the form of relic abundances of helium and other light nuclei. The nucleosynthesis is really a very narrow bottleneck for all cosmological models, and one which has amply confirmed the standard Big Bang model. We find that the baryonic matter present since nucleosynthesis is completely insufficient to close the Universe.

In Section 6.5 we try to understand baryosynthesis and the absence of antimatter.

6.1 Planck Time

Recall that in quantum mechanics it is always possible to associate the mass of a particle M with a wave having the *Compton wavelength*

$$\lambda = \frac{\hbar}{Mc}. \quad (6.1)$$

In other words, for a particle of mass M , quantum effects become important at distances of the order of λ . On the other hand, gravitational effects are important at distances of the order of the Schwarzschild radius. Equating the two distances, we find the scale at which quantum effects and gravitational effects are of equal importance. This defines the *Planck mass*

$$M_{\text{p}} = \sqrt{\hbar c/G} = 1.221 \times 10^{19} \text{ GeV } c^{-2}. \quad (6.2)$$

From this we can derive the Planck energy $M_{\text{p}}c^2$ and the *Planck time*

$$t_{\text{p}} = \lambda_{\text{p}}/c = 5.31 \times 10^{-44} \text{ s}. \quad (6.3)$$

Later on we shall make frequent use of quantities at the Planck scale. The reason for associating these scales with Planck's name is that he was the first to notice that the combination of fundamental constants

$$\lambda_{\text{p}} = \sqrt{\hbar G/c^3} = 1.62 \times 10^{-35} \text{ m} \quad (6.4)$$

yielded a natural length scale. The particle symmetry at Planck time is characterized by all fields except the inflaton field being exactly massless. Only when this symmetry is spontaneously broken in the transition to a lower temperature phase do some particles become massive.

Unfortunately, there is as yet no theory including quantum mechanics and gravitation. Thus we have no description of the Universe before the Planck time, nor of the Big Bang, because of a lack of theoretical tools.

There is strong evidence that this cosmic inflation happened when the energy scale of the Universe was about three orders of magnitude lower than the Planck scale. We shall postpone the discussion of inflation to Chapter 7.

6.2 The Primordial Hot Plasma

The primeval Universe may have developed through phases when some symmetry was exact, followed by other phases when that symmetry was broken. The early cosmology would then be described by a sequence of *phase transitions*. Symmetry breaking may occur through a *first-order phase transition*, in which the field tunnels through a potential barrier, or through a *second-order phase transition*, in which the field evolves smoothly from one state to another, following the curve of the potential.

An important bookkeeping parameter at all times is the temperature, T . When we follow the history of the Universe as a function of T , we are following a trajectory in

space-time which may be passing through regions of different vacua. In the simple model of symmetry breaking by a real scalar field ϕ the T -dependence may be put in explicitly, as well as other dependencies (denoted by dots),

$$V(\phi, T, \text{etc.}) = -\frac{1}{2}\mu^2\phi^2 + \frac{1}{4}\lambda\phi^4 + \frac{1}{8}\lambda T^2\phi^2 + \dots \quad (6.5)$$

As time decreases T increases, the vacuum expectation value ϕ_0 decreases, so that finally, in the early Universe, the true minimum of the potential is the trivial one at $\phi = 0$. This occurs above a *critical temperature* of

$$T_c = 2\mu/\sqrt{\lambda}. \quad (6.6)$$

An example of this behavior is illustrated by the potentials in Figure 6.1. The different curves correspond to different temperatures. At the highest temperature, the only minimum is at $\phi = 0$ but, as the temperature decreases, a new minimum develops spontaneously. If there is more than one minimum, only one of these is stable. A classical example of an unstable minimum is a steam bubble in boiling water.

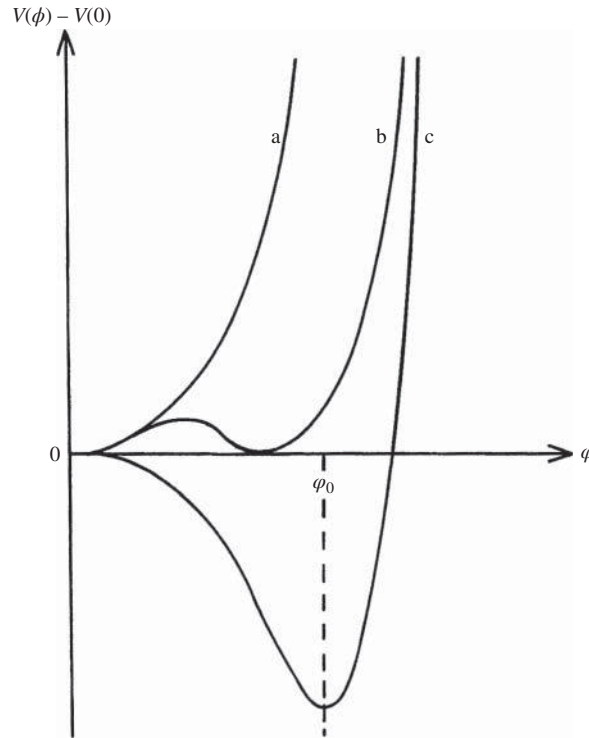


Figure 6.1 Effective scalar potentials. The different curves correspond to different temperatures. At the highest temperature (a) the only minimum is at $\phi = 0$ but, as the temperature decreases (b), a new minimum develops spontaneously. Finally, in (c), a stable minimum is reached at ϕ_0 .

Particle Physics Symmetries. Above $E \approx 10^{15}$ GeV the Universe was filled with a plasma of a *Grand Unified field* of extreme heat and pressure, and occupying an exceedingly small volume. Neither elementary particles nor their more elementary constituents were stable under such conditions, only quantum packets of radiative energy incessantly exchanging energy existed. Only when the Universe had cooled to below 1 TeV the leptons characterized by their electroweak interactions and the quarks characterized by their color interactions appeared as constituents in the plasma.

In the energy range 10^{12} GeV $\lesssim E \lesssim 10^{16}$ GeV the strengths of the electro-weak force and the quark color force experienced by photons, leptons and quarks, might have been unified in a unique symmetry group G_{GUT} , with one coupling constant g . All the leptons and quarks should be components of that same field.

If there is such a symmetry group its mathematical form is not known. Although all GUTs are designed to answer some of the questions and relate some of the parameters in the standard model, they still have the drawback of introducing large numbers of new particles, vector bosons and Higgs scalars, all of which have to be discovered.

The global symmetry group in particle physics may have been of the form

$$G_s \equiv SU(3)_c \otimes SU(2)_w \otimes U(1)_{B-L}. \quad (6.7)$$

This is referred to as the *standard model*, where the subscript c stands for *color*, w for *weak interactions* and $B - L$ for *Baryon and Lepton number*, respectively. It is considered to be an exact symmetry in the $1 \text{ TeV} \lesssim E \lesssim 10^{11-12} \text{ TeV}$ energy range. It is laboratory physics that has led us to construct the standard model, which is fairly well understood although experimental information above 1 TeV is lacking. Moreover, the standard model leaves a number of questions open which one would very much like to have answered within a GUT model.

There are too many free parameters in the standard model. Why are the electric charges such as they are? How many Higgs scalars are there, more than the one discovered in 2012? Why do the leptons and the quarks come in three families. The GUT symmetry group would require only one family, but if nature provides us with more, why are there precisely three? What is the reason for CP violation? The only hint is that CP violation seems to require (but not explain) at least three families of quarks.

The big question is what new physics may appear in the enormous range between 1 TeV and the GUT energy 10^{14} GeV. The possibility that nothing new appears is called ‘the desert’.

The new physics could be a higher symmetry which would be broken at the lower end of this energy range. Somewhere there would then be a phase transition between the exactly symmetric phase and the spontaneously broken phase. Even in the case of a ‘desert’, one expects a phase transition at 10^{14} or 10^{15} GeV. One effect which finds its explanation in processes at about 10^{14} GeV is the remarkable absence of antimatter in the Universe.

If the GUT symmetry breaks down to the standard model through intermediate steps, the phenomenology could be very rich. For instance, there are *subconstituent models* building leptons and quarks out of elementary particles of one level deeper elementarity. These subconstituents would freeze out at some intermediate energy,

condensing into lepton and quark bound states. The forces binding them are in some models called *technicolor forces*.

Below the energy $E \approx 1$ TeV we encounter the phase transition between exact and spontaneously broken $SU(2)_w \otimes U(1)_{B-L}$ symmetry. The end of electroweak unification is marked by the massification of the vector boson fields, the scalar Higgs fields and the fermion fields.

One much discussed extension to the standard model is *supersymmetry* (SUSY). This brings in a large number of new particles, some of which should be seen in this temperature range. In this theory there is a conserved multiplicative quantum number, *R parity*, defined by

$$R = (-1)^{3B+L+2s}, \quad (6.8)$$

where B , L and s are baryon number, lepton number and spin, respectively. All known particles have $R = +1$, but the theory allows an equal number of supersymmetric partners, *sparticles*, having $R = -1$. Conservation of R ensures that SUSY sparticles can only be produced pairwise, as sparticle–antisparticle pairs. The lightest sparticles must therefore be stable, just as the lightest particles are.

One motivation for introducing this intermediate scale is the *hierarchy problem*: why is m_p so enormously much larger than m_w ? And why is V_{Coulomb} so much larger than V_{Newton} ? SUSY has so many free parameters that it can ‘naturally’ explain these problems.

Thermal Equilibrium. Motion of particles under electromagnetic interaction is described by the Maxwell–Lorentz equations. The motion of a particle in a central field of force F , as for instance an electron of charge e moving at a distance r around an almost static proton, is approximated well by the *Coulomb force*

$$F = \frac{e^2}{r^2}. \quad (6.9)$$

Note that this has the same form as Newton’s law of gravitation, Equation (1.28). In the electromagnetic case the strength of the interaction is e^2 , whereas the strength of the gravitational interaction is $Gm_m c$. These two *coupling constants* are expressed in completely different units because they apply to systems of completely different sizes. For the physics of radiation, the gravitational interaction can be completely neglected but, for the dynamics of the expansion of the Universe, only the gravitational interaction is important because celestial objects are electrically neutral.

The photons and the first particles which formed momentarily in the primordial plasma were incessantly colliding and exchanging energy and momentum at relativistic speeds. A few collisions were sufficient to distribute the available energy evenly among them. Their reaction rates were much greater than the Hubble expansion rate, so thermal equilibrium should have been maintained in any local comoving volume element dV .

There was no net inflow or outflow of energy, which defines the expansion as *adiabatic* *adiabatic*, as was done in Equation (5.25). The law of conservation of energy

[Equation (5.24)], also called the *first law of thermodynamics*, followed by assuming that matter behaved as an expanding nonviscous fluid at constant pressure p .

When the collisions resulted in a stable energy spectrum, *thermal equilibrium* was established and the photons had the *blackbody spectrum* derived in 1900 by Max Planck.

Let the number of photons of energy $h\nu$ per unit volume and frequency interval be $n_\gamma(\nu)$. Then the photon number density in the frequency interval $(\nu, \nu + d\nu)$ is

$$n_\gamma(\nu) d\nu = \frac{8\pi}{c^3} \frac{\nu^2 d\nu}{e^{h\nu/kT} - 1}. \quad (6.10)$$

At the end of the nineteenth century some 40 years was spent trying to find this formula using trial and error. With the benefit of hindsight, the derivation is straightforward, based on classical thermodynamics as well as on quantum mechanics, unknown at Planck's time.

Note that Planck's formula depends on only one parameter, the temperature T . Thus the energy spectrum of photons in thermal equilibrium is completely characterized by its temperature T . The distribution [Equation (6.10)] peaks at the frequency

$$\nu_{\max} \simeq 3.32 \times 10^{10} T \quad (6.11)$$

in units of Hertz or cycles per second, where T is given in degrees Kelvin.

The total number of photons per unit volume, or the *number density* N_γ , is found by integrating this spectrum over all frequencies:

$$N_\gamma = \int_0^\infty n_\gamma(\nu) d\nu \simeq 1.202 \frac{2}{\pi^2} \left(\frac{kT}{c\hbar} \right)^3. \quad (6.12)$$

Here \hbar represents the reduced Planck's constant $\hbar = h/2\pi$. The temperature T may be converted into units of energy by the dimensional relation $E = kT$. The solution of the integral in this equation can be given in terms of Riemann's zeta-function; $\zeta(3) \approx 1.2020$.

Since each photon of frequency ν is a quantum of energy $h\nu$ (this is the interpretation Planck was led to, much to his own dismay, because it was in obvious conflict with classical ideas of energy as a continuously distributed quantity), the total energy density of radiation is given by the *Stefan-Boltzmann law* after *Josef Stefan* (1835–1893) and *Ludwig Boltzmann* (1844–1906),

$$\epsilon_r = \int_0^\infty h\nu n_\gamma(\nu) d\nu = \frac{\pi^2}{15} \frac{k^4 T^4}{\hbar^3 c^3} \equiv a_S T^4, \quad (6.13)$$

where all the constants are lumped into Stefan's constant

$$a_S = 4723 \text{ eV m}^{-3} \text{ K}^{-4}.$$

The blackbody spectrum is shown in Figure 8.1.

If the expansion is adiabatic and the pressure p is constant so that $d(pV) = p dV$, we recover the *second law of thermodynamics*.

The second law of thermodynamics states in particular that entropy cannot decrease in a closed system. The particles in a plasma possess maximum entropy when thermal equilibrium has been established. The assumption that the Universe expands

adiabatically is certainly very good during the radiation-dominated era when the fluid was composed of photons and elementary particles in thermal equilibrium. However, the Universe, black holes and for instance voids are generalized thermodynamic systems for which the equations of thermodynamics must be written in a complete form involving a factor due to the expansion of the Universe and a redefinition of entropy.

This is also true during the matter-dominated era before matter clouds start to contract into galaxies under the influence of gravity. Even on a very large scale we may consider the galaxies forming a homogeneous ‘fluid’, an idealization as good as the cosmological principle that forms the basis of all our discussions. In fact, we have already relied on this assumption in the derivation of the Einstein equation and in the discussion of equations of state. However, the pressure in the ‘fluid’ of galaxies of density N is negligibly small, because it is caused by their random motion, just as the pressure in a gas is due to the random motion of the molecules. Since the average peculiar velocities $\langle v \rangle$ of the galaxies are of the order of $10^{-3}c$, the ratio of pressure $p = m\langle v \rangle^2 N$ to matter density ρ gives an equation of state (Problem 5) of the order of

$$w \approx \frac{m\langle v \rangle^2 N}{\rho c^2} = \frac{\langle v \rangle^2}{c^2} \approx 10^{-6}.$$

We have already relied on this value in the case of a matter-dominated universe when deriving Equation (5.30).

Energy Density. Let us compare the energy densities of radiation and matter. The energy density of electromagnetic radiation corresponding to one photon in a volume V is

$$\rho_r c^2 \equiv \varepsilon_r = \frac{h\nu}{V} = \frac{hc}{V\lambda}. \quad (6.14)$$

In an expanding universe with cosmic scale factor a , all distances scale as a and so does the wavelength λ . The volume V then scales as a^3 ; thus ε_r scales as a^{-4} . Here and in the following the subscript ‘r’ stands for radiation and relativistic particles, while ‘m’ stands for nonrelativistic (cold) matter.

Statistical mechanics tells us that the pressure in a nonviscous fluid is related to the energy density by the equation of state [Equation (5.32)]

$$p = \frac{1}{3}\varepsilon, \quad (6.15)$$

where the factor $\frac{1}{3}$ comes from averaging over the three spatial directions. Thus pressure also scales as a^{-4} , so that it will become even more negligible in the future than it is now. The energy density of matter,

$$\rho_m c^2 = \frac{mc^2}{V}, \quad (6.16)$$

also decreases with time, but only with the power a^{-3} . Thus the ratio of radiation energy to matter scales as a^{-1} :

$$\frac{\varepsilon_r}{\rho_m} \propto \frac{a^{-4}}{a^{-3}} \propto a^{-1}. \quad (6.17)$$

The present radiation energy density is predominantly in the form of microwaves and infrared light. The change from radiation domination to matter domination is gradual: at $t = 1000$ yr the radiation fraction was about 90%, at $t = 2$ Myr only about 10% (see Figure 7.3).

Relativistic versus Nonrelativistic Particles. It is important to distinguish between relativistic and nonrelativistic particles because their energy spectra in thermal equilibrium are different. A coarse rule is that a particle is nonrelativistic when its kinetic energy is small in comparison with its mass, and relativistic when $E \gtrsim 10mc^2$. The masses of some cosmologically important particles are given in Table A.4. For comparison, the equivalent temperatures are also given. This gives a rough idea of the temperature of the heat bath when the respective particle is nonrelativistic.

The adiabaticity condition $\dot{S} = 0$ can be applied to both relativistic and nonrelativistic particles. Let us first consider the relativistic particles which dominate the radiation era. Recall from Equation (3.9) that the energy of a particle depends on two terms, mass and kinetic energy,

$$E = \sqrt{m^2c^4 + P^2c^2}, \quad (6.18)$$

where P is momentum. For massless particles such as the photons, the mass term is of course absent; for relativistic particles it can be neglected.

When p is a constant we can write the law of energy conservation

$$dE = -pdV. \quad (6.19)$$

Replacing E by the energy density ϵ_r times the volume $a^3 \equiv V$, the above law becomes

$$d(a^3\epsilon_r) = -p da^3. \quad (6.20)$$

Substituting ϵ_r for the pressure p from the equation of state Equation (6.15) we obtain

$$a^3 d\epsilon_r + \epsilon_r da^3 = -\frac{1}{3}\epsilon_r da^3,$$

or

$$\frac{d\epsilon_r}{\epsilon_r} = -\frac{4}{3} \frac{da^3}{a^3}. \quad (6.21)$$

The solution to this equation is

$$\epsilon_r \propto a^{-4}, \quad (6.22)$$

in agreement with our previous finding. We have in fact already used this result in Equation (5.33).

For nonrelativistic particles the situation is different. Their kinetic energy ϵ_{kin} is small, so that the mass term in Equation (6.18) can no longer be neglected. The motion of n particles per unit volume is then characterized by a temperature T_m , causing a pressure

$$p = nkT_m. \quad (6.23)$$

Note that T_m is not the temperature of matter in thermal equilibrium, but rather a bookkeeping device needed for dimensional reasons. The equation of state differs from that of radiation and relativistic matter, Equation (6.15), by a factor of 2:

$$p = \frac{2}{3} \epsilon_{\text{kin}}.$$

Including the mass term of the n particles, the energy density of nonrelativistic matter becomes

$$\rho_m \equiv \epsilon_m = nmc^2 + \frac{3}{2}nkT_m. \quad (6.24)$$

Substituting Equations (6.23) and (6.24) into Equation (6.20) we obtain

$$d(a^3 nmc^2) + \frac{3}{2}d(a^3 nkT_m) = -nkT_m da^3. \quad (6.25)$$

Let us assume that the total number of particles always remains the same: in a scattering reaction there are then always two particles coming in, and two going out, whatever their types. This is not strictly true because there also exist other types of reactions producing more than two particles in the final state. However, let us assume that the total number of particles in the volume V under consideration is $N = Vn$, and that N is constant during the adiabatic expansion,

$$dN = d(Vn) = \frac{4}{3}\pi d(a^3 n) = 0. \quad (6.26)$$

The first term in Equation (6.25) then vanishes and we are left with

$$\frac{3}{2}a^3 dT_m = -T_m d(a^3),$$

or

$$\frac{3}{2} \frac{dT_m}{T_m} = -\frac{d(a^3)}{a^3}.$$

The solution to this differential equation is of the form

$$T_m \propto a^{-2}. \quad (6.27)$$

Thus we see that the temperature of nonrelativistic matter has a different dependence on the scale of expansion than does the temperature of radiation. This has profound implications for one of the most serious problems in thermodynamics in the nineteenth century.

The number density of relativistic particles other than photons is given by distributions very similar to the Planck distribution. Let us replace the photon energy $h\nu$ in Equation (6.10) by E , which is given by the relativistic expression in Equation (6.18). Noting that the kinematic variable is now the three-momentum $p = |\mathbf{p}|$ (since for relativistic particles we can ignore the mass), we can replace Planck's distribution by the number density of particle species i with momentum between p and $p + dp$,

$$n_i(p) dp = \frac{8\pi}{h^3} \frac{n_{\text{spin},i}}{2} \frac{p^2 dp}{e^{E_i(p)/kT_i} \pm 1}. \quad (6.28)$$

The \pm sign is ‘-’ for bosons and ‘+’ for fermions, and the name for these distributions are the *Bose distribution* and the *Fermi distribution*, respectively. The Fermi distribution in the above form is actually a special case: it holds when the number of charged fermions equals the number of corresponding neutral fermions (the ‘chemical potentials’ vanish). In the following we shall need only that case.

The number density N of nonrelativistic particles of mass m is given by the *Maxwell–Boltzmann* distribution for an ideal, nondegenerate gas. Starting from Equation (6.28) we note that for nonrelativistic particles the energy kT is smaller than the mass, so that the term ± 1 in can be neglected in comparison with the exponential. Rewriting the Fermi distribution as a function of temperature rather than of momentum we obtain the Maxwell–Boltzmann distribution

$$N = n_{\text{spin}} \frac{(2\pi mkT)^{3/2}}{(hc)^3} e^{-E_i/kT}. \quad (6.29)$$

Note that because of the exponential term the number density falls exponentially as temperature falls. *James Clerk Maxwell* (1831–1879) was a contemporary of Stefan and Boltzmann.

Thermal Death. Suppose that the Universe starts out at some time with γ rays at high energy and electrons at rest. This would be a highly ordered nonequilibrium system. The photons would obviously quickly distribute some of their energy to the electrons via various scattering interactions. Thus the original order would decrease, and the randomness or disorder would increase. The second law of thermodynamics states that any isolated system left by itself can only change towards greater disorder. The measure of disorder is entropy; thus the law says that entropy cannot decrease.

The counterexample which living organisms seem to furnish, since they build up ordered systems, is not valid. This is because no living organism exists in isolation; it consumes nutrients and produces waste. Thus, establishing that a living organism indeed increases entropy would require measurement of a much larger system, certainly not smaller than the Solar System.

It now seems to follow from the second law of thermodynamics that all energy would ultimately distribute itself evenly throughout the Universe, so that no further temperature differences would exist. The discoverer of the law of conservation of energy, *Hermann von Helmholtz* (1821–1894), came to the distressing conclusion in 1854 that ‘from this point on, the Universe will be falling into a state of eternal rest’. This state was named *thermal death*, and it preoccupied greatly both philosophers and scientists during the nineteenth century.

Now we see that this pessimistic conclusion was premature. Since E scales as a^{-1} it follows that also the temperature of radiation T_r scales as a^{-1} . Thus, from the time when the temperatures of matter and radiation were equal,

$$T_m = T_r,$$

we see from Equation (6.27) that the adiabatic expansion of the Universe causes matter to cool faster than radiation. Thus cold matter and hot radiation in an expanding

Universe are not and will never be in thermal equilibrium on a cosmic timescale. This result permits us to solve the adiabatic equations of cold matter and hot radiation separately, as we in fact have.

6.3 Electroweak Interactions

In *quantum electrodynamics* (QED) the electromagnetic field is mediated by photons which are emitted by a charged particle and absorbed very shortly afterwards by another. Photons with such a brief existence during an interaction are called *virtual*, in contrast to real photons.

Virtual particles do not travel freely to or from the interaction region. Energy is not conserved in the production of virtual particles. This is possible because the energy imbalance arising at the creation of the virtual particle is compensated for when it is annihilated, so that the real particles emerging from the interaction region possess the same amount of energy as those entering the region. We have already met this argument in the discussion of Hawking radiation from black holes.

However, nature impedes the creation of very huge energy imbalances. For example, the masses of the *vector bosons* W^\pm and Z^0 mediating the electroweak interactions are almost 100 GeV. Reactions at much lower energies involving virtual vector bosons are therefore severely impeded, and much less frequent than electromagnetic interactions. For this reason such interactions are called *weak interactions*.

Real photons interact only with charged particles such as protons p , electrons e^- and their oppositely charged *antiparticles*, the *anti-proton* \bar{p} and the *positron* e^+ . An example is the elastic *Compton scattering* of electrons by photons:

$$\gamma + e^\pm \rightarrow \gamma + e^\pm. \quad (6.30)$$

As a result of virtual intermediate states neutral particles may exhibit electromagnetic properties such as magnetic moment. When an electron is captured by a free proton, they form a bound state, a hydrogen atom which is a very stable system. An electron and a positron may also form a bound atom-like state called *positronium*. This is a very unstable system: the electron and positron are antiparticles, so they rapidly end up annihilating each other according to the reaction

$$e^- + e^+ \rightarrow \gamma + \gamma. \quad (6.31)$$

Since the total energy is conserved, the annihilation results in two (or three) photons possessing all the energy and flying away with it at the speed of light.

The reverse reaction is also possible. A photon may convert briefly into a virtual e^-e^+ pair, and another photon may collide with either one of these charged particles, knocking them out of the virtual state, thus creating a free electron–positron pair:

$$\gamma + \gamma \rightarrow e^- + e^+. \quad (6.32)$$

This requires the energy of each photon to equal at least the electron (positron) mass, 0.51 MeV. If the photon energy is in excess of 0.51 MeV the e^-e^+ pair will not be created at rest, but both particles will acquire kinetic energy.

Protons and anti-protons have electromagnetic interactions similar to positrons and electrons. They can also annihilate into photons, or for instance into an electron-positron pair via the mediation of a virtual photon,

$$p + \bar{p} \rightarrow \gamma_{\text{virtual}} \rightarrow e^- + e^+. \quad (6.33)$$

The reverse reaction is also possible, provided the electron and positron possess enough kinetic energy to create a proton, or 938.3 MeV.

Conservation laws. Note that the total electric charge is conserved throughout the reactions in Equations (6.29)–(6.32). Electric charge can never disappear nor arise out of neutral vacuum, but it can easily move from a charged particle to a neutral one as long as that does not violate the conservation of total charge in the reaction. In the annihilation of an e^-e^+ pair into photons, all charges do indeed vanish, but only because the sum of the charges was zero to start with.

There are two further conservation laws governing the behavior of baryons and leptons.

- (i) B or *baryon number* is conserved. This forbids the total number of baryons minus anti-baryons from changing in particle reactions. To help the bookkeeping in particle reactions one assigns the value $B = 1$ to baryons and $B = -1$ to anti-baryons in a way analogous to the assignment of electric charges. Photons and leptons have $B = 0$.
- (ii) L_l or *l-lepton number* is conserved for each of the lepton flavours $l = e, \mu, \tau$. This forbids the total number of l -leptons minus \bar{l} -anti-leptons from changing in particle reactions. We assign $L_e = 1$ to e^- and ν_e , $L_e = -1$ to e^+ and $\bar{\nu}_e$, and correspondingly to the members of the μ and τ families. Photons and baryons have no lepton numbers.

However, there is an amendment to this rule, caused by the complications in the physics of neutrinos. Although the flavour state l is conserved in neutrino reactions, it is not conserved in free flight. To observe the flavour state l of neutrinos is not the same as observing the neutrino mass states. There are three neutrino mass states called ν_1, ν_2, ν_3 , which are not identical to the flavour states; rather, they are quantum-mechanical superpositions of them. The states are entangled in such a way that a pure mass state is a mixture of flavour states, and vice versa. Roughly, the ν_μ is the mixture of $\frac{1}{4}\nu_1, \frac{1}{4}\nu_2$ and $\frac{1}{2}\nu_3$.

All leptons participate in the weak interactions mediated by the heavy virtual vector bosons W^\pm and Z^0 and the scalar *Higgs boson* H^0 . The Z^0 is just like a photon except that it is very massive, about 91 GeV, and the two W^\pm are its 10 GeV lighter charged partners. The mass of the H^0 is 125 GeV. All these bosons freeze out of thermal equilibrium at $E \approx 100$ GeV.

There is no difference between weak and electromagnetic interactions: there are charged-current electroweak interactions mediated by the W^\pm , and neutral-current interactions mediated by the Z^0 , the H^0 and the γ . However, the electroweak symmetry is imperfect because of the very different masses.

Weak leptonic reactions are

$$e^\pm + \bar{\nu}_e^{(-)} \rightarrow e^\pm + \nu_e^{(-)}, \quad (6.34)$$

$$\bar{\nu}_e^{(-)} + \nu_e^{(-)} \rightarrow \bar{\nu}_e^{(-)} + \nu_e^{(-)}, \quad (6.35)$$

where $\bar{\nu}_e^{(-)}$ stands for $\bar{\nu}_e$ or $\bar{\nu}_e$. There is also the annihilation reaction

$$e^- + e^+ \rightarrow \nu_e + \bar{\nu}_e, \quad (6.36)$$

as well as the reverse pair production reaction.

Similar reactions apply to the two other lepton families, replacing e above by μ or τ , respectively. Note that the ν_e can scatter against electrons by the W^\pm exchange and Z^0 exchange, respectively. In contrast, ν_μ and ν_τ can only scatter by the Z^0 exchange diagram, because of the separate conservation of lepton-family numbers.

Spin and Statistics. The leptons and nucleons all have two spin states each. In the following we shall refer to them as *fermions*, after *Enrico Fermi* (1901–1954), whereas the photon, the *Higgs boson* and the W and Z are *bosons*, after *Satyendranath Bose* (1894–1974).

The difference between bosons and fermions is deep and fundamental. The number of spin states is even for fermions, odd for bosons (except the photon). They behave differently in a statistical ensemble. Fermions have antiparticles which most bosons do not. The *fermion number* is conserved, indeed separately for leptons and baryons. The number of bosons is not conserved; for instance, in pp collisions one can produce any number of pions and photons.

Two identical fermions refuse to get close to one another. This is the *Pauli exclusion force* for which Wolfgang Pauli (1900–1958) received the Nobel prize, and which is responsible for the electron *degeneracy pressure* in white dwarfs and the neutron degeneracy pressure in neutron stars. A gas of free electrons will exhibit pressure even at a temperature of absolute zero. According to quantum mechanics, particles never have exactly zero velocity: they always carry out random motions, causing pressure. For electrons in a high-density medium such as a white dwarf with density $10^6 \rho_\odot$, the degeneracy pressure is much larger than the thermal pressure, and it is enough to balance the pressure of gravity.

Bosons do not feel such a force, nothing inhibits them getting close to each other. The massive vector bosons W^\pm and Z^0 have three spin or polarization states: the *transversal* (vertical and horizontal) states which the photons also have, and the *longitudinal state* along the direction of motion, which the photon is lacking.

The number of distinct states or *degrees of freedom*, g , of photons in a statistical ensemble (in a plasma, say) is two. In general, due to the intricacies of quantum statistics, the degrees of freedom are the product of the number of spin states, n_{spin} , the number n_{anti} which is 2 for particles with distinct antiparticles and otherwise 1, and a factor $n_{\text{Pauli}} = \frac{7}{8}$, which only enters for fermions obeying Fermi–Dirac statistics. For bosons this factor is unity. The degrees of freedom are then

$$g = n_{\text{spin}} n_{\text{anti}} n_{\text{Pauli}}, \quad (6.37)$$

tabulated in the fifth column of Table A.5.

For each species of relativistic fermions participating in the thermal equilibrium there is a specific number density. To find the total number density of particles sharing the available energy we have to count each particle species i weighted by the corresponding degrees of freedom g_i .

Equation (6.13) with a factor g_i explicitly visible is

$$\varepsilon_i = \frac{1}{2} g_i a_S T^4, \quad (6.38)$$

where a_S is Stefan's constant. It turns out that this expression gives the correct energy density for every particle species if we insert its respective value of g_i from Table A.5.

Equation (6.12) can be correspondingly generalized to relativistic fermions. Their number density is

$$N_f = \frac{3}{4} N_\gamma. \quad (6.39)$$

In general, the primordial plasma was a mixture of particles, of which some are relativistic and some nonrelativistic at a given temperature. Since the number density of a nonrelativistic particle [given by the Maxwell–Boltzmann distribution, Equation (6.28)] is exponentially smaller than that of a relativistic particle, it is a good approximation to ignore nonrelativistic particles. Different species i with mass m_i have a number density which depends on m_i/T , and they may have a thermal distribution with a temperature T_i different from that of the photons. Let us define the *effective degrees of freedom* of the mixture as

$$g_* = \sum_{\text{bosons } i} g_i + \sum_{\text{fermions } j} g_j \left(\frac{T_j}{T} \right)^4. \quad (6.40)$$

As explained in the context of Equation (6.37) the sum over fermions includes a factor $\frac{7}{8}$, accounting for the difference between Fermi and Bose statistics. The factor $(T_j/T)^4$ applies only to neutrinos, which obtain a different temperature from the photons when they freeze out from the plasma (as we shall see later). Thus the energy density of the radiation in the plasma is

$$\varepsilon_r = \frac{1}{2} g_* a_S T^4. \quad (6.41)$$

The sum of degrees of freedom of a system of particles is of course the number of particles multiplied by the degrees of freedom per particle. Independently of the law of conservation of energy, the conservation of entropy implies that the energy is distributed equally between all degrees of freedom present in such a way that a change in degrees of freedom is accompanied by a change in random motion, or equivalently in temperature.

Thus entropy is related to order: the more degrees of freedom there are present, the more randomness or disorder the system possesses. When an assembly of particles (such as the molecules in a gas) does not possess energy other than kinetic energy (heat), its entropy is maximal when thermal equilibrium is reached. For a system of gravitating bodies, entropy increases by clumping, maximal entropy corresponding to a black hole.

Temperature and Time Scales. Let us now derive a relation between the temperature scale and the timescale. We have already found the relation (5.40) between the size scale R and the timescale t during the radiation era,

$$a(t) \propto \sqrt{t}, \quad (6.42)$$

where we choose to omit the proportionality factor. The Hubble parameter can then be written

$$H = \frac{\dot{a}}{a} = \frac{1}{2t}. \quad (6.43)$$

Note that the proportionality factor omitted in Equation (6.42) has dropped out.

In Equation (5.35) we noted that the curvature term kc^2/R^2 in Friedmann's equations is negligibly small at early times during the radiation era. We then obtained the dynamical relation

$$\frac{\dot{a}}{a} = \left(\frac{8\pi G}{3} \rho \right)^{1/2}. \quad (6.44)$$

Inserting Equation (6.43) on the left and replacing the energy density ρ on the right by ϵ_r/c^2 , we find the relation sought between photon temperature and time:

$$\frac{1}{t} = \sqrt{\frac{16\pi G a_S}{3c^2} g_*} T^2 = 3.07 \times 10^{-21} \sqrt{g_*} \frac{T^2}{[\text{K}^2]} [\text{s}^{-1}]. \quad (6.45)$$

The sum of degrees of freedom of a system of particles is of course the number of particles multiplied by the degrees of freedom per particle. Independently of the law of conservation of energy, the conservation of entropy implies that the energy is distributed equally between all degrees of freedom present in such a way that a change in degrees of freedom is accompanied by a change in random motion, or equivalently in temperature.

Thus entropy is related to order: the more degrees of freedom there are present, the more randomness or disorder the system possesses. When an assembly of particles (such as the molecules in a gas) does not possess energy other than kinetic energy (heat), its entropy is maximal when thermal equilibrium is reached. For a system of gravitating bodies, entropy increases by clumping, maximal entropy corresponding to a black hole.

Cooling Plasma. At a temperature of 10^{11} K, which corresponds to a mean energy of about 300 MeV, the particles contributing to the effective degrees of freedom g_* are the photon, three charged leptons, three neutrinos (not counting their three inert right-handed components), the six quarks with three colors each, the gluon of eight colors and two spin states, the scalar Higgs boson H^0 , and the vector gauge bosons W^\pm, Z^0 . Thus the effective degrees of freedom is given by

$$g_* = 2 + 3 \times \frac{7}{2} + 3 \times \frac{7}{4} + 6 \times 3 \times \frac{7}{2} + 8 \times 2 + 1 + 3 \times 3 = 106.75. \quad (6.46)$$

Above the QCD-hadron phase transition at 200 MeV the hadrons are represented by their free quark subconstituents which contribute more degrees of freedom than

hadrons, forming a dense medium of *quark matter*. Below the transition the quarks become bound in hadronic matter; the separation between quarks in the nucleons increases, and the interaction between any two quarks in a nucleon cease to interact with quarks in neighboring nucleons. The degrees of freedom $g_*(T)$ are then reduced steeply to $69/4$ as is shown in Figure 6.2 [1].

Quark matter may still exist today in the core of cold but dense stellar objects such as neutron stars.

The color symmetry $SU(3)_c$ is valid and unbroken at this temperature, but it is in no way apparent, because the hadrons are color-neutral singlets. The color force mediated by gluons is also not apparent: a vestige of QCD from earlier epochs remains in the form of strong interactions between hadrons. It appears that this force is mediated by mesons, themselves quark bound states.

There is no trace of the weak symmetry $SU(2)_w$, so the weak and electromagnetic interactions look quite different. Their strengths are very different, and the masses of the leptons are very different. Only the electromagnetic gauge symmetry $U(1)$ is exactly valid, as is testified to by the conservation of electric charge.

All electrons and photons have an energy below the threshold for proton-anti-proton production [see Equation (6.33)]. Then the number of protons, neutrons and anti-nucleons will no longer increase as a result of thermal collisions, they can only decrease.

We can follow the evolution of the function $g_*(T)$ in Figure 6.2 [1]. Actually all the particles in thermal equilibrium contribute, not only those accounted for in Equation (6.61), but also heavier particles which are thermalized by energetic photons and pions in the tail of their respective Boltzmann distributions.

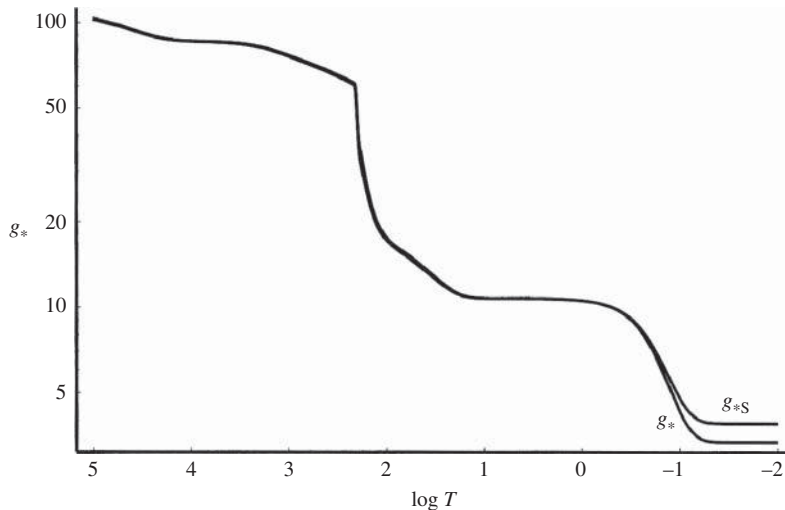


Figure 6.2 The evolution of the effective degrees of freedom contributing to the energy density, $g_*(T)$ and to the entropy density, $g_{*S}(T)$, as functions of $\log T$, where the temperature is in units of MeV [1].

At this time the number density of nucleons decreases quickly because they have become nonrelativistic. Consequently, they have a larger probability of annihilating into lepton pairs, pion pairs or photons. Their number density is then no longer given by the Fermi distribution [Equation (6.28)], but by the Maxwell–Boltzmann distribution, Equation (6.29). As can be seen from the latter, when T drops below the mass, the number density decreases rapidly because of the exponential factor. If there had been exactly the same number of nucleons and anti-nucleons, we would not expect many nucleons to remain to form matter. But, since we live in a matter-dominated Universe, there must have been some excess of nucleons early on. Note that neutrons and protons exist in equal numbers at the time under consideration.

Although the nucleons are very few, they still participate in electromagnetic reactions such as the elastic scattering of electrons,

$$e^\pm + p \rightarrow e^\pm + p, \quad (6.47)$$

and in weak *charged current* reactions in which charged leptons and nucleons change into their neutral partners, and vice versa, as in

$$e^- + p \rightarrow \nu_e + n, \quad (6.48)$$

$$\bar{\nu}_e + p \rightarrow e^+ + n. \quad (6.49)$$

Other such reactions are obtained by reversing the arrows, and by replacing e^\pm by μ^\pm or ν_e by ν_μ or ν_τ . The nucleons still participate in thermal equilibrium, but they are too few to play any role in the thermal history any more.

Below the pion mass (actually at about 70 MeV) the temperature in the Universe cools below the threshold for pion production:

$$(e^- + e^+) \text{ or } (\mu^- + \mu^+) \rightarrow \gamma_{\text{virtual}} \rightarrow \pi^+ + \pi^-. \quad (6.50)$$

The reversed reactions, pion annihilation, still operate, reducing the number of pions. However, they disappear even faster by decay. This is always the fate when such lighter states are available, energy and momentum, as well as quantum numbers such as electric charge, baryon number and lepton numbers, being conserved. The pion, the muon and the tau lepton are examples of this. The pion decays mainly by the reactions

$$\pi^- \rightarrow \mu^- + \bar{\nu}_\mu, \quad \pi^+ \rightarrow \mu^+ + \nu_\mu. \quad (6.51)$$

Thus g_* decreases by 3 to $\frac{57}{4}$. The difference in mass between the initial pion and the final state particles is

$$m_\pi - m_\mu - m_\nu = (139.6 - 105.7 - 0.0) \text{ MeV} = 33.9 \text{ MeV}, \quad (6.52)$$

so 33.9 MeV is available as kinetic energy to the muon and the neutrino. This makes it very easy for the π^\pm to decay, and in consequence its mean life is short, only 0.026 μs (the π^0 decays even faster). This is much less than the age of the Universe at 140 MeV, which is 23 μs from Equation (6.45). Note that the appearance of a charged lepton in the final state forces the simultaneous appearance of its anti-neutrino in order to conserve lepton number.

Also, the muons decay fast compared with the age of the Universe, with a lifetime of $2.2 \mu\text{s}$, by the processes

$$\mu^- \rightarrow e^- + \bar{\nu}_e + \nu_\mu, \quad \mu^+ \rightarrow e^+ + \nu_e + \bar{\nu}_\mu. \quad (6.53)$$

Almost the entire mass of the muon, or 105.7 MeV , is available as kinetic energy to the final state particles. This is the reason for its short mean life. Here again the conservation of lepton numbers, separately for the e -family and the μ -family, is observed.

Below the muon mass (actually at about 50 MeV), the temperature in the Universe cools below the threshold for muon-pair production:

$$e^- + e^+ \rightarrow \gamma_{\text{virtual}} \rightarrow \mu^+ + \mu^-. \quad (6.54)$$

The time elapsed is less than a millisecond. When the muons have disappeared, we can reduce g_* by $\frac{7}{2}$ to $\frac{43}{4}$.

From the reactions in Equations (6.51) and (6.53) we see that the end products of pion and muon decay are stable electrons and neutrinos. The lightest neutrino ν_1 is certainly stable, and the same is probably also true for ν_2 and ν_3 . When this has taken place we are left with those neutrinos and electrons that only participate in weak reactions, with photons and with a very small number of nucleons. The number density of each lepton species is about the same as that of photons.

6.4 Photon and Lepton Decoupling

The considerations about which particles participate in thermal equilibrium at a given time depend on two timescales: the *reaction rate* of the particle, taking into account the reactions which are possible at that energy, and the *expansion rate* of the Universe. If the reaction rate is slow compared with the expansion rate, the distance between particles grows so fast that they cannot find each other.

Reaction Rates. The expansion rate is given by $H = \dot{a}/a$, and its temperature dependence by Equations (6.44) and (6.45). The average reaction rate can be written

$$\Gamma = \langle Nv\sigma(E) \rangle, \quad (6.55)$$

where $\sigma(E)$ is the reaction cross-section. The product of $\sigma(E)$ and the velocity v of the particle varies over the thermal distribution, so one has to average over it, as is indicated by the angle brackets. Multiplying this product by the number density N of particles per m^3 , one obtains the mean rate Γ of reacting particles per second, or the mean collision time between collisions, Γ^{-1} .

The weak interaction cross-section turns out to be proportional to T^2 ,

$$\sigma \simeq \frac{G_F^2 (kT)^2}{\pi (\hbar c)^4}, \quad (6.56)$$

where G_F is the *Fermi coupling* measuring the strength of the weak interaction. The number density of the neutrinos is proportional to T^3 according to Equations (6.12) and (6.39). The reaction rate of neutrinos of all flavours then falls with decreasing temperature as T^5 .

The condition for a given species of particle to remain in thermal equilibrium is then that the reaction rate Γ is larger than the expansion rate H , or equivalently that Γ^{-1} does not exceed the Hubble distance H^{-1} ,

$$\frac{\Gamma}{H} \gtrsim 1. \quad (6.57)$$

Inserting the T^5 dependence of the weak interaction rate Γ_{wi} and the T^2 dependence of the expansion rate H from Equation (6.45), we obtain

$$\frac{\Gamma_{\text{wi}}}{H} \propto T^3. \quad (6.58)$$

Thus there may be a temperature small enough that the condition in Equation (5.62) is no longer fulfilled.

Photon Reheating. Photons with energies below the electron mass can no longer produce e^+e^- pairs, but the energy exchange between photons and electrons still continues by Compton scattering, reaction in Equation (6.30), or *Thomson scattering*, as it is called at very low energies. Electromagnetic cross-sections (subscript ‘em’) are proportional to T^{-2} , and the reaction rate is then proportional to T , so

$$\frac{\Gamma_{\text{em}}}{H} \propto \frac{1}{T}.$$

Contrary to the weak interaction case in Equation (6.58), the condition in Equation (6.57) is then satisfied for all temperatures, so electromagnetic interactions never freeze out. Electrons only decouple when they form neutral atoms during the *Recombination Era* and cease to scatter photons. The term recombination is slightly misleading, because the electrons have never been combined into atoms before. The term comes from laboratory physics, where free electrons and *ionized* atoms are created by heating matter (and upon subsequent cooling the electrons and ions recombine into atoms) or from so-called HII regions, where interstellar plasma is ionized by ultraviolet radiation and characteristic *recombination radiation* is emitted when electrons and ions re-form.

The exothermic electron–positron annihilation, reaction in Equation (6.31), is now of mounting importance, creating new photons with energy 0.51 MeV. This is higher than the ambient photon temperature at that time, so the photon population gets reheated. To see just how important this reheating is, let us turn to the law of conservation of entropy.

Making use of the equation of state for relativistic particles (6.15), the entropy can be written

$$S = \frac{4V}{3kT} \epsilon_{\text{plasma}}.$$

Substituting the expression for ϵ_{plasma} from Equation (6.41) one obtains

$$S = \frac{2g_*}{3} \frac{V a_S T^4}{kT}, \quad (6.59)$$

which is valid where we can ignore nonrelativistic particles. Now $a_S T^4$ is the energy density, so $V a_S T^4$ is energy, just like kT , and thus $V a_S T^4/kT$ is a constant. g_* is also

a constant, except at the thresholds where particle species decouple. The physical meaning of entropy of a system is really its degrees of freedom multiplied by some constant, as one sees here. In Equation (6.12) we saw that the entropy density can also be written

$$s \equiv \frac{S}{V} = \frac{3}{2} \zeta(3) g_\gamma N_\gamma, \quad (6.60)$$

where N_γ is the number density of photons. Between two decoupling thresholds we then have

$$\frac{dS}{dt} = \frac{d}{dt} \left(\frac{2g_*}{3} \frac{V a_S T^3}{k} \right) = 0. \quad (6.61)$$

The second law of thermodynamics requires that entropy should be conserved in reversible processes, also at thresholds where g_* changes. This is only possible if T also changes in such a way that $g_* T^3$ remains a constant. When a relativistic particle becomes nonrelativistic and disappears, its entropy is shared between the particles remaining in thermal contact, causing some slight slowdown in the cooling rate. Photons never become nonrelativistic; neither do the practically massless neutrinos, and therefore they continue to share the entropy of the Universe, each species conserving its entropy separately.

Let us now apply this argument to the situation when the positrons and most of the electrons disappear by annihilation below 0.2 MeV. We denote temperatures and entropies just above this energy by a subscript '+', and below it by '-'. Above this energy, the particles in thermal equilibrium are γ , e^- , e^+ . Then the entropy

$$S = \frac{2}{3} \left(2 + \frac{7}{2} \right) \frac{V a_S T_+^3}{k}. \quad (6.62)$$

Below that energy, only photons contribute the factor $g_* = 2$. Consequently, the ratio of entropies S_+ and S_- is

$$\frac{S_+}{S_-} = \frac{11}{4} \left(\frac{T_+}{T_-} \right)^3. \quad (6.63)$$

But entropy must be conserved so this ratio must be unity. It then follows that

$$T_- = \left(\frac{11}{4} \right)^{1/3} T_+ = 1.40 T_+. \quad (6.64)$$

Thus the temperature T_γ of the photons increases by a factor of 1.40 as the Universe cools below the threshold for electron–positron pair production. Actually, the temperature increase is so small and so gradual that it only slows down the cooling rate temporarily.

Neutrino Decoupling. When the neutrinos no longer obey the condition in Equation (6.57) they *decouple* or *freeze out* from all interactions, and begin a free expansion. The decoupling of ν_μ and ν_τ occurs at about 3.8 MeV, whereas the ν_e decouple at 2.3 MeV. This can be depicted as a set of connecting baths containing different particles, and having valves which close at given temperatures (see Figure 6.3).

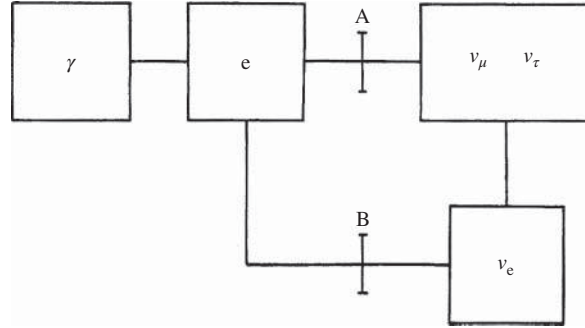


Figure 6.3 A system of communicating vessels illustrating particles in thermal equilibrium (from K. Kainulainen, unpublished research). At about 3.8 MeV, valve A closes so that ν_μ and ν_τ decouple. At 2.3 MeV, valve B closes so that ν_e also decouples, leaving only e^- and γ in thermal contact.

At decoupling, the neutrinos are still relativistic, since they are so light (Table A.3). Thus their energy distribution is given by the Fermi distribution, Equation (6.28), and their temperature equals that of the photons, $T_\nu = T_\gamma$, decreasing with the increasing scale of the Universe as a^{-1} . But the neutrinos do not participate in the reheating process, and they do not share the entropy of the photons, so from now on they remain colder than the photons:

$$T_\nu = T_\gamma/1.40. \quad (6.65)$$

The number density N_ν of neutrinos can be calculated as in Equation (6.12) using Equation (6.10), except that the -1 term in Equation (6.10) has to be replaced by $+1$, which is required for fermions [see Equation (6.28)]. In the number density distributions [Equations (6.10) and (6.28)], we have ignored possible chemical potentials for all fermions, which one can do for a thermal radiation background; for neutrinos it is an unproven assumption that nonetheless appears in reasonable agreement with their oscillation parameters.

The result is that N_ν is a factor of $\frac{3}{4}$ times N_γ at the same temperature. Taking the difference between temperatures T_ν and T_γ into account and noting from Equation (6.12) that N_γ is proportional to T^3 , one finds

$$N_\nu = \frac{3}{4} \frac{4}{11} N_\gamma. \quad (6.66)$$

After decoupling, the neutrino contribution to g_* decreases because the ratio T_i/T in Equation (6.40) is now less than one. Thus the present value is

$$g_*(T_0) = 2 + 3 \frac{7}{4} \left(\frac{4}{11} \right)^{4/3} = 3.36. \quad (6.67)$$

The entropy density also depends on g_* , but now the temperature dependence for the neutrino contribution in Equation (6.40) is $(T_i/T)^3$ rather than a power of four. The effective degrees of freedom are in that case given by Equation (6.67) if the power $\frac{4}{3}$ is replaced by 1. This curve is denoted g_{*S} in Figure 6.2.

The density parameter is

$$\Omega_\nu = \frac{3}{11} \frac{N_\gamma}{\rho_c h^2} \sum_i m_i. \quad (6.68)$$

Recombination Era. As long as there are free electrons, the primordial photons are thermalized by Thomson scattering against them, and this prohibits the electrons from decoupling, in contrast to neutrinos. Each scattering polarizes the photons, but on average this is washed out. The electromagnetic reaction rate is much higher than the weak reaction rate of the neutrinos; in fact, it is higher than the expansion rate of the Universe, so the condition in Equation (6.57) is fulfilled.

Eventually the Universe expands and cools to such an extent, to about 1000 K, that electrons are captured into atomic orbits, primarily by protons but also by the trace amounts of ionized helium and other light nuclei. This process is referred to as recombination. Unlike the unstable particles n , π , μ that decay spontaneously liberating kinetic energy in exothermic reactions, the hydrogen atom H is a *bound state* of a proton and an electron. Its mass is less than the p and e^- masses together,

$$m_H - m_p - m_e = -13.59 \text{ eV}, \quad (6.69)$$

so it cannot disintegrate spontaneously into a free proton and a free electron. The mass difference in Equation (6.69) is the *binding energy* of the hydrogen atom.

The physics of recombination is somewhat subtle. Initially one might think that recombination occurs when the photon temperature drops below 13.59 eV, making formation of neutral hydrogen energetically favorable. Two characteristics of the physics push the recombination temperature lower. The first, and the easiest to elaborate, is that there are vastly more photons than electrons and so in thermal equilibrium even a small proportion of high-energy photons are sufficient to maintain effectively complete ionization. Photons in thermal equilibrium have the blackbody spectrum given by Equation (6.10). Even for photon temperatures somewhat below 13.59 eV there will be enough highly energetic photons in the Wein tail (as the high-energy section is termed) to ionize any neutral hydrogen. The large amount of entropy in the Universe also favors free protons and electrons.

With respect to the thermal history of the Universe, the fact that photons do not scatter against neutral atoms is critical. As recombination proceeds and the number of electrons falls, matter and radiation decouple. This has two results. First, with matter and radiation no longer in thermal equilibrium the thermal history of the two proceed independently. Perturbations in matter are no longer damped by interaction with radiation and any such perturbations can grow into structures through gravitational instability. Decoupling thus initiates the period of structure formation that has led to our present Universe being populated with stars, galaxies, galaxy clusters and so on.

The second result is that, with photons no longer scattering against a sea of electrons, the photons can stream freely through the Universe; upon recombination, the Universe becomes transparent to light. Prior to recombination, the Universe was opaque to electromagnetic radiation (although not to neutrinos) and it would have been impossible to do astronomy if this situation had persisted until today. The freely

streaming photons from this era form the CMB radiation and their point of last contact with matter forms an isotropic *last scattering surface* (LSS). The era of recombination provides a crucial observational limit beyond which we cannot hope to see using electromagnetic radiation. The LSS is not a sharp boundary and does not exist at a unique redshift. Once the conditions for recombination have been met in one portion they should be also met in any other—thus isotropicity. The reason for seeing this as a theoretical spherical shell is that only one specific shell at one specific distance can be seen by any specific observer at any specific time. The photons from the LSS preserve the polarization they incurred in the last Thomson scattering. This remaining primordial polarization is an interesting detectable signal, albeit much weaker than the intensity of the thermalized radiation.

The LSS of the Universe has an exact analogue in the surface of the Sun. Photons inside the Sun are continuously scattered, so it takes millions of years for some photons to reach the surface. But once they do not scatter any more they continue in straight lines (really on geodesics) towards us. Therefore, we can see the surface of the Sun, which is the LSS of the solar photons, but we cannot see the solar interior. We can also observe that sunlight is linearly polarized. In contrast, neutrinos hardly scatter at all in the Sun, thus neutrino radiation brings us a clear (albeit faint with present neutrino detectors) picture of the interior of the Sun.

Equilibrium Theory. An analysis based on thermal equilibrium, the *Saha equation*, implies that the temperature must fall to about 0.3 eV before the proportion of high-energy photons falls sufficiently to allow recombination to occur. The Saha analysis also implies that the time (or energy or redshift) for decoupling and last scattering depend on cosmological parameters such as the total cosmic density parameter Ω_0 , the baryon density Ω_b , and the Hubble parameter. However, a second feature of the physics of recombination implies that the equilibrium analysis itself is not sufficient.

The Saha analysis describes the initial phase of departure from complete ionization but, as recombination proceeds, the assumption of equilibrium ceases to be appropriate [see, e.g., Equation (6.9)]. Paradoxically, the problem is that electromagnetic interactions are too fast (in contrast with the weak interaction that freezes out from equilibrium because of a small cross-section). A single recombination directly to the ground state would produce a photon with energy greater than the 13.59 eV binding energy and this photon would travel until it encountered a neutral atom and ionized it. This implies that recombination in an infinite static universe would have to proceed by smaller intermediate steps (thus not directly to the ground state).

In fact the situation is even worse, because reaching the ground state by single photon emission requires transition from the 2P to 1S levels and thus production of photons with energy at least 10.2 eV (Lyman α with $\lambda = 1216 \text{ \AA}$). As these photons become abundant they will re-ionize any neutral hydrogen through multiple absorption and so it would seem that recombination will be, at a minimum, severely impeded. (Recombination in a finite HII region is different because the Ly α photons can escape.)

There is an alternative path, however. Two-photon emission generated by the $2S \rightarrow 1S$ transition produces lower-energy photons. The process is slow (with a lifetime of approximately 0.1 s), so recombination proceeds at a rate quite different from the

Saha prediction. Consequently, all the times predicted by this nonequilibrium analysis differs notably from the Saha prediction, but, interestingly, in such a way that the times of decoupling and last scattering have practically no dependence on cosmological parameters.

Summary. Recombination, decoupling and last scattering do not occur at the exactly same time. It should also be noted that these terms are often used interchangeably in the literature, so what we refer to as the LSS is sometimes called the time of recombination or decoupling. Recombination can be defined as the time when 90% of the electrons have combined into neutral atoms. This occurred at redshift

$$z_{\text{recombination}} \approx 1100. \quad (6.70)$$

Last scattering is defined as the time when photons start to stream freely. This occurred at redshift

$$z_{\text{LSS}} = 1089, \quad (6.71)$$

when the Universe was $379\,000(\Omega_0 h^2)^{-1/2}$ years old and at a temperature of 0.26 eV, thus right after the recombination time.

Decoupling is defined as the time when the reaction rate (scattering) fell below the expansion rate of the Universe and matter fell out of thermal equilibrium with photons. This occurred at redshift

$$z_{\text{decoupling}} \approx 890. \quad (6.72)$$

All three events depend on the number of free electrons (the ionization fraction) but in slightly different ways. As a result these events do not occur at exactly the same time, but close enough to explain the confusion in naming.

The build up of structures started right after decoupling. Reionization of the Universe occurred at $z_{\text{reionization}} = 20 \pm 5$ and the matter structures were in place at $z_{\text{structures}} = 5$.

6.5 Big Bang Nucleosynthesis

Let us now turn to the fate of the remaining nucleons. Note that the charged current reactions in Equations (6.47) and (6.48) changing a proton to a neutron are *endothermic*: they require some input energy to provide for the mass difference. In the reaction in Equation (6.48) this difference is 0.8 MeV and in reaction in Equation (6.49) it is 1.8 MeV (use the masses in Table A.4!). The reversed reactions are *exothermic*. They liberate energy and they can then always proceed without any threshold limitation.

The neutrons and protons are then nonrelativistic, so their number densities are each given by Maxwell–Boltzmann distributions [Equation (6.29)]. Their ratio in equilibrium is given by

$$\frac{N_n}{N_p} = \left(\frac{m_n}{m_p} \right)^{3/2} \exp \left(-\frac{m_n - m_p}{kT} \right). \quad (6.73)$$

At energies of the order of $m_n - m_p = 1.293$ MeV or less, this ratio is dominated by the exponential. Thus, at $kT = 0.8$ MeV, the ratio has dropped from 1 to $\frac{1}{5}$. As the Universe cools and the energy approaches 0.8 MeV, the endothermic neutron-producing reactions stop, one by one. Then no more neutrons are produced but some of those that already exist get converted into protons in the exothermic reactions.

Nuclear Fusion. Already, at a few MeV, nuclear *fusion reactions* start to build up light elements. These reactions are exothermic: when a neutron and a proton fuse into a bound state some of the nucleonic matter is converted into pure energy according to Einstein's formula [Equation (3.8)]. This binding energy of the *deuteron* d,

$$m_p + m_n - m_d = 2.22 \text{ MeV}, \quad (6.74)$$

is liberated in the form of radiation:



The deuteron is also written ${}^2\text{H}^+$ in general nuclear physics notation, where the superscript $A = 2$ indicates the number of nucleons and the electric charge is given by the superscript '+'. The bound state formed by a deuteron and an electron is the *deuterium* atom ${}^2\text{H}$, which of course is electrically neutral. Although the deuterons are formed in very small quantities, they are of crucial importance to the final composition of matter.

As long as photons of 2.22 MeV or more are available, the reaction in Equation (6.75) can go the other way: the deuterons *photodisintegrate* into free protons and neutrons. Even when the mean temperature of radiation drops considerably below 2.22 MeV, there is still a high-energy tail of the Planck distribution containing hard γ rays which destroy the deuterons as fast as they are produced.

All evidence suggests that the number density of baryons, or equivalently nucleons, is today very small. In particular, we are able to calculate it to within a factor $\Omega_B h^2$,

$$N_B = \frac{\rho_B}{m_B} = \frac{\Omega_B \rho_c}{m_B} \simeq 11.3 \Omega_B h^2 \text{ m}^{-3}. \quad (6.76)$$

At the end of this section we shall discuss the value of the baryon density parameter Ω_B , which is a small percentage.

The photon number density today is $N_\gamma = 4.11 \times 10^8$ per m^3 from Equation (6.12). It is clear then that N_B/N_γ is such a small figure that only an extremely tiny fraction of the high-energy tail of the photon distribution may contain sufficiently many hard γ rays to photodisintegrate the deuterons. However, the 2.22 MeV photons created in photodisintegration do not thermalize, so they will continue to photodisintegrate deuterium until they have been redshifted below this threshold. Another obstacle to permanent deuteron production is the high entropy per nucleon in the Universe. Each time a deuteron is produced, the degrees of freedom decrease, and so the entropy must be shared among the remaining nucleons. This raises their temperature, counteracting the formation of deuterons. Detailed calculations show that deuteron production becomes thermodynamically favored only at 0.07 MeV. Thus, although deuterons are

avored on energetic grounds already at 2 MeV, free nucleons continue to be favored by the high entropy down to 0.07 MeV.

Other nuclear fusion reactions also commence at a few MeV. The npp bound state ${}^3\text{He}^{++}$ is produced in the fusion of two deuterons,



where the final-state particles share the binding energy

$$2m_p + m_n - m({}^3\text{He}^{++}) = 7.72 \text{ MeV}. \quad (6.79)$$

This reaction is also hampered by the large entropy per nucleon, so it becomes thermodynamically favored only at 0.11 MeV.

The nnp bound state ${}^3\text{H}^+$, or *triton* t, is the ionized *tritium* atom, ${}^3\text{H}$. It is produced in the fusion reactions



with the binding energy

$$m_p + 2m_n - m_t = 8.48 \text{ MeV}. \quad (6.83)$$

A very stable nucleus is the nnpp bound state ${}^4\text{He}^{++}$ with a very large binding energy,

$$2m_p + 2m_n - m({}^4\text{He}^{++}) = 28.3 \text{ MeV}. \quad (6.84)$$

Once its production is favored by the entropy law, at about 0.28 MeV, there are no more γ rays left that are hard enough to photodisintegrate it. From the examples set by the deuteron fusion reactions above, it may seem that ${}^4\text{He}^{++}$ would be most naturally produced in the reaction



However, ${}^3\text{He}^{++}$ and ${}^3\text{H}^+$ production is preferred over deuteron fusion, so ${}^4\text{He}^{++}$ is only produced in a second step when these nuclei become abundant. The reactions are then



The delay before these reactions start is often referred to as the *deuterium bottleneck*.

Below 0.8 MeV occasional weak interactions in the high-energy tails of the lepton and nucleon Fermi distributions reduce the n/p ratio further, but no longer by

the exponential factor in Equation (6.73). The neutrons also decay into protons by *beta decay*,

$$n \rightarrow e^- + \bar{\nu}_e + p, \tag{6.90}$$

liberating

$$m_n - m_p - m_e - m_{\nu} = 0.8 \text{ MeV} \tag{6.91}$$

of kinetic energy in the process. This amount is very small compared with the neutron mass of 939.6 MeV. In consequence the decay is inhibited and very slow: the neutron mean life is 887 s. In comparison with the age of the Universe, which at this time is a few tens of seconds, the neutrons are essentially stable. The protons are stable even on scales of billions of years, so their number is not going to decrease by decay.

At 0.1 MeV, when the temperature is 1.2×10^9 K and the time elapsed since the Big Bang is a little over two minutes, the beta decays have reduced the neutron/proton ratio to its final value:

$$\frac{N_n}{N_p} \simeq \frac{1}{7}. \tag{6.92}$$

The temperature dependence of this ratio, as well as the equilibrium (Maxwell-Boltzmann) ratio, is shown in Figure 6.4.

These remaining neutrons have no time to decay before they fuse into deuterons and subsequently into $^4\text{He}^{++}$. There they stayed until today because bound neutrons do not decay. The same number of protons as neutrons go into ^4He , and the remaining free protons are the nuclei of future hydrogen atoms. Thus the end result of the nucleosynthesis taking place between 100 and 700 s after the Big Bang is a Universe composed almost entirely of hydrogen and helium ions. But why not heavier nuclei?

It is an odd circumstance of nature that, although there exist stable nuclei composed of $A = 1, 2, 3$ and 4 nucleons, no nucleus of $A = 5$, nor really of $A = 8$, exists. In between these gaps, there exist the stable nuclei ^6Li and ^7Li and also ^7Be which

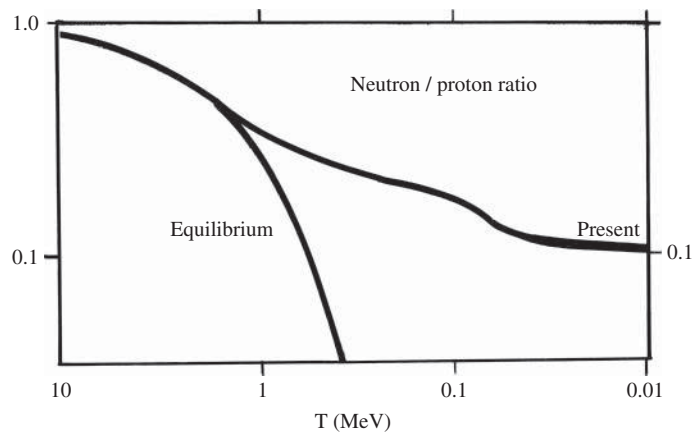


Figure 6.4 The *neutron/proton* ratio from 10 Mev to the present ≈ 10 keV. The equilibrium n/p ratio is also shown.

was stable in comparison with the age of the universe. At $A = 8$ the nuclei ${}^8\text{B}$ and ${}^8\text{Li}$ have lifetimes near one second which are not quite total bottlenecks. Because of these gaps and bottlenecks and because ${}^4\text{He}$ is so strongly bound, nucleosynthesis essentially stops after ${}^4\text{He}$ production. Only minute quantities of the stable nuclei ${}^2\text{H}$, ${}^3\text{He}$ and ${}^7\text{Li}$ remain with important relic abundances until today.

The fusion rate at energy E of two nuclei of charges Z_1, Z_2 is proportional to the *Gamow penetration factor*

$$\exp\left(-\frac{2Z_1Z_2}{\sqrt{E}}\right). \quad (6.93)$$

Thus as the energy decreases, the fusion of nuclei other than the very lightest ones becomes rapidly improbable.

Relic ${}^4\text{He}$ Abundance. The relic abundances of the light elements bear an important testimony of the n/p ratio at the time of the nucleosynthesis when the Universe was only a few minutes old. In fact, this is the earliest testimony of the Big Bang we have. Recombination occurred some 300 000 years later, when the stable ions captured all the electrons to become neutral atoms. The CMB testimony is from that time. There is also more recent information available in galaxy cluster observations from $z < 0.25$.

From the ratio in Equation (6.92) we obtain immediately the ratio of ${}^4\text{He}$ to ${}^1\text{H}$:

$$X_4 \equiv \frac{N({}^4\text{He})}{N({}^1\text{H})} = \frac{N_n/2}{N_p - N_n} \simeq \frac{1}{12}. \quad (6.94)$$

The number of ${}^4\text{He}$ nuclei is clearly half the number of neutrons when the minute amounts of ${}^2\text{H}$, ${}^3\text{He}$ and ${}^7\text{Li}$ are neglected. The same number of protons as neutrons go into ${}^4\text{He}$, thus the excess number of protons becoming hydrogen is $N_p - N_n$. The ratio of mass in ${}^4\text{He}$ to total mass in ${}^1\text{H}$ and ${}^4\text{He}$ is

$$Y_4 \equiv \frac{4X_4}{1 + 4X_4} = 0.2477 \pm 0.0001. \quad (6.95)$$

from CMB measurements [8]. This is in some tension with the value 0.2551 ± 0.0022 from [9]. This is a function of the ratio of baryons to photons

$$\eta \equiv \frac{N_b}{N_\gamma} \simeq 2.75 \times 10^{-8} \Omega_b h^2, \quad (6.96)$$

using N_γ from Table A.6.

The helium mass abundance Y_4 depends sensitively on several parameters. If the number of baryons increases, Ω_b and η also increase, and the entropy per baryon decreases. Since the large entropy per baryon was the main obstacle to early deuteron and helium production, the consequence is that helium production can start earlier. But then the neutrons would have had less time to β -decay, so the neutron/proton ratio would be larger than $\frac{1}{7}$. It follows that more helium will be produced: Y_4 increases.

The abundances of the light elements also depend on the neutron mean life τ_n and on the number of neutrino families F_ν , both of which were poorly known until 1990. Although τ_n is now known to 1‰ [4], and F_ν is known to be $3 \pm 4\%$ [2], it may be

instructive to follow the arguments about how they affect the value of Y_4 . Let us rewrite the decoupling condition [Equation (6.58)] for neutrons

$$\frac{\Gamma_{\text{wi}}}{H} = AT_{\text{d}}^3, \quad (6.97)$$

where A is the proportionality constant left out of Equation (5.64) and T_{d} is the decoupling temperature. An increase in the neutron mean life implies a decrease in the reaction rate Γ_{wi} and therefore a decrease in A . At temperature T_{d} the ratio of the reaction rate to the expansion rate is unity; thus

$$T_{\text{d}} = A^{-1/3}. \quad (6.98)$$

Hence a longer neutron mean life implies a higher decoupling temperature and an earlier decoupling time. As we have already seen, an earlier start of helium production leads to an increase in Y_4 .

The expansion rate H of the Universe is, according to Equations (6.43) and (6.45), proportional to $\sqrt{g_*}$, which in turn depends on the number of neutrino families F_ν . In Equations (6.46) we had set $F_\nu = 3$. Thus, if there were more than three neutrino families, H would increase and A would decrease with the same consequences as in the previous example. Similarly, if the number of neutrinos were very different from the number of anti-neutrinos, contrary to the assumptions in standard Big Bang cosmology, H would also increase.

Light Element Abundance Observations. The value of $\Omega_{\text{b}}h^2$ (or η) is obtained in direct measurements of the relic abundances of ^4He , ^3He , ^2H or D, and ^7Li from the time when the Universe was only a few minutes old. Although the ^4He mass ratio Y_4 is about 0.25, the ^3He and ^2H mass ratios are less than 10^{-4} and the ^7Li mass ratio as small as a few times 10^{-10} , they all agree remarkably well on a common value for η .

If the observed abundances are indeed of cosmological origin, they must not be affected significantly by later stellar processes. The helium isotopes ^3He and ^4He cannot be destroyed easily but they are continuously produced in stellar interiors. Some recent helium is blown off from supernova progenitors, but that fraction can be corrected for by observing the total abundance in hydrogen clouds of different ages and extrapolating to time zero. The remainder is then primordial helium emanating from BBN. On the other hand, the deuterium abundance can only decrease; it is easily burned to ^3He in later stellar events. Measurements of the deuterium abundance are quite uncertain because of systematic errors. The case of ^7Li is complicated because some fraction is due to later galactic cosmic ray spallation products.

The ^4He abundance is easiest to observe, but it is also least sensitive to $\Omega_{\text{b}}h^2$, its dependence is logarithmic, so only very precise measurements are relevant. The best 'laboratories' for measuring the ^4He abundance are a class of low-luminosity dwarf galaxies called blue compact dwarf (BCD) galaxies, which undergo an intense burst of star formation in a very compact region. The BCDs are among the most metal-deficient gas-rich galaxies known (astronomers call all elements heavier than helium *metals*). Since their gas has not been processed through many generations of stars, it should approximate well the pristine primordial gas.

The ^3He isotope can be seen in galactic star-forming regions containing ionized hydrogen (HII), in the local interstellar medium and in planetary nebulae. Because HII regions are objects of zero age when compared with the age of the Galaxy, their elemental abundances can be considered typical of primordial conditions.

The ^7Li isotope is observed at the surface of the oldest stars. Since the age of stars can be judged by the presence of metals, the constancy of this isotope has been interpreted as being representative of the primordial abundance.

The strongest constraint on the baryonic density comes from the primordial deuterium abundance. Ultraviolet light with a continuous flat spectrum emitted from objects at distances of $z \approx 2 - 3.5$ will be seen redshifted into the red range of the visible spectrum. Photoelectric absorption in intervening hydrogen along the line of sight then causes a sharp cut-off at $\lambda = 91.2$ nm, the *Lyman limit*. This can be used to select objects of a given type, which indeed are star-forming galaxies. Deuterium is observed as a Lyman- α feature in the absorption spectra of high-redshift quasars. A recent analysis [5] gives

$$\Omega_b(^2\text{H})h^2 = 0.020 \pm 0.001, \quad (6.99)$$

which is more precise than any other determination. The information from the other light nucleids are in good agreement. The values of η and Ω_b in Table A.6 come from a combined fit to ^2H data, CMB and large-scale structure. We defer that discussion to Section 8.4.

In Figure 6.5 the history of the Universe is summarized in nomograms relating the scales of temperature, energy, size, density and time. Note that so far we have only covered the events which occurred between 10^{11} K and 10^3 K.

Nuclear synthesis also goes on inside stars where the gravitational contraction increases the pressure and temperature so that the fusion process does not stop with helium. Our Sun is burning hydrogen to helium, which lasts about 10^{10} yr, a time span which is very dependent on the mass of the star. After that, helium burns to carbon in typically 10^6 yr, carbon to oxygen and neon in 10^4 yr, those to silicon in 10 yr, and silicon to iron in 10 h, whereafter the fusion chain stops. The heavier elements have to be synthesized much later in supernova explosions, and all elements heavier than lithium have to be distributed into the intergalactic medium within the first billion years.

To sum up, Big Bang cosmology makes some very important predictions. The Universe today should still be filled with freely streaming primordial photon (and neutrino) radiation with a blackbody spectrum [Equation (6.10)] of temperature related to the age of the Universe and a polarization correlated to the temperature. This relic CMB radiation (as well as the relic neutrino radiation) should be essentially isotropic since it originated in the now spherical shell of the LSS. In particular, it should be uncorrelated to the radiation from foreground sources of later date, such as our Galaxy. We shall later see that these predictions have been verified for the photons (but not yet for the neutrinos).

A very important conclusion from BBN is that the Universe contains surprisingly little baryonic matter! Either the Universe is then indeed open, or there must exist other types of nonbaryonic, gravitating matter.

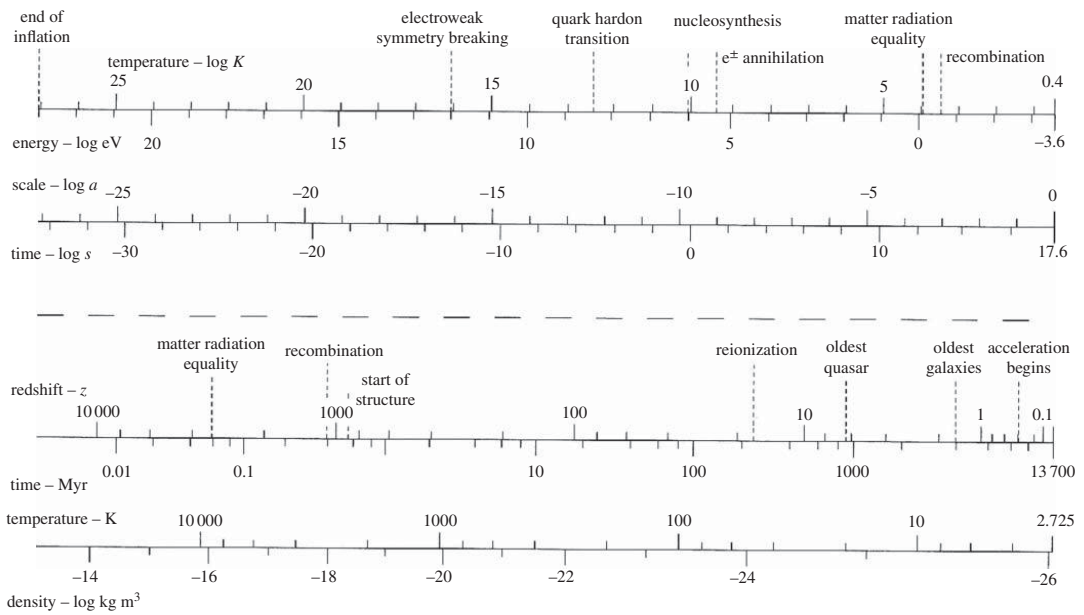


Figure 6.5 The 'complete' history of the Universe. The scales above the dashed line cover the time from the end of inflation to the present, those below from the end of radiation dominance to the present. Input constants (not completely updated): $H_0 = 0.71 \text{ km s}^{-1} \text{ Mpc}^{-1}$, $\Omega_\gamma = 0.73$, $T_0 = 2.275 \text{ K}$, $t_0 = 13.7 \text{ Gyr}$, and g_* , g_{*S} from Figure 6.2.

6.6 Baryosynthesis and Antimatter Generation

In Equation (6.96) we noted that the ratio η of the baryon number density N_B to the photon number density N_γ is very small. We shall see later that it is even two orders of magnitude smaller.

The Ratio of Baryons to Photons. Before the baryons were formed (at about 200 MeV), the conserved baryon number B was carried by quarks. Thus the total value of B carried by all protons and neutrons today should equal the total value carried by all quarks. It is very surprising then that η should be so small today, because when the quarks, leptons and photons were in thermal equilibrium there should have been equal numbers of quarks and anti-quarks, leptons and anti-leptons, and the value of B was equal to the total leptonic number L ,

$$N_B = N_{\bar{B}} = N_L \approx N_\gamma \quad (6.100)$$

because of the $U(1)_{B-L}$ symmetry.

When the baryons and anti-baryons became nonrelativistic, the numbers of baryons and anti-baryons were reduced by annihilation, so N_B decreased rapidly by the exponential factor in the Maxwell–Boltzmann distribution [Equation (6.29)]. The number density of photons N_γ is given by Equation (6.12) at all temperatures. Thus the temperature dependence of η should be

$$\eta = \frac{N_B}{N_\gamma} = \frac{\sqrt{2\pi}}{4.808} \left(\frac{m_N}{kT} \right)^{3/2} e^{-m_N/kT}. \quad (6.101)$$

When the annihilation rate became slower than the expansion rate, the value of η was frozen, and thus comparable to its value today. The freeze-out occurs at about 20 MeV, when η has reached the value

$$\eta \simeq 6.8 \times 10^{-19}. \quad (6.102)$$

But this is a factor 9×10^8 too small! Thus something must be seriously wrong with our initial condition [Equation (6.100)].

Baryon–Anti-baryon Asymmetry. The other surprising thing is that no anti-baryons seem to have survived. At temperatures above 200 MeV, quarks and anti-quarks were in thermal equilibrium with the photons because of reactions such as Equations (6.30)–(6.32), as well as

$$\gamma + \gamma \leftrightarrow q + \bar{q}. \quad (6.103)$$

These reactions conserve baryon number, so every quark produced or annihilated is accompanied by one anti-quark produced or annihilated. Since all quarks and anti-quarks did not have time to annihilate one another, it would be reasonable to expect equal number densities of baryons or anti-baryons today.

But we know that the Earth is only matter and not antimatter. The solar wind does not produce annihilations with the Earth or with the other planets, so we know that

the Solar System is matter. Since no gamma rays are produced in the interaction of the solar wind with the local interstellar medium, we know that the interstellar medium, and hence our Galaxy, is matter. The main evidence that other galaxies are not composed of antimatter comes from cosmic rays. Our Galaxy contains free protons (cosmic rays) with a known velocity spectrum very near the speed of light. A fraction of these particles have sufficiently high velocities to escape from the Galaxy. These protons would annihilate with antimatter cosmic rays in the intergalactic medium or in collisions with antimatter galaxies if they existed, and would produce characteristic gamma rays many orders of magnitude more frequently than have been seen. The observed ratio is

$$\frac{N_{\bar{B}}}{N_B} = 10^{-5} - 10^{-4},$$

depending on the kinetic energy. This small number is fully consistent with all the observed anti-protons having been produced by energetic cosmic rays in the Earth's atmosphere, and it essentially rules out the possibility that other galaxies emit cosmic rays composed of antimatter. There are many other pieces of evidence against antimatter, but the above arguments are the strongest.

We are then faced with two big questions. What caused the large value of η ? And why does the Universe not contain antimatter, anti-quarks and positrons. The only reasonable conclusion is that $N_{\bar{B}}$ and N_B must have started out slightly different while they were in thermal equilibrium, by the amount

$$N_B - N_{\bar{B}} \simeq \eta N_\gamma. \quad (6.104)$$

Subsequently most anti-baryons were annihilated, and the small excess ηN_γ of baryons is what remained. This idea is fine, but the basic problem has not been removed; we have only pushed it further back to earlier times, to some early $B\bar{B}$ -asymmetric phase transition.

Primeval Asymmetry Generation. Let us consider theories in which a $B\bar{B}$ -asymmetry could arise. For this three conditions must be met.

First, the theory must contain reactions violating baryon number conservation. Grand unified theories are obvious candidates for a reason we have already met in Section 6.2. We noted there that GUTs are symmetric with respect to leptons and quarks, because they are components of the same field and GUT forces do not see any difference. Consequently, GUTs contain leptoquarks X, Y which transform quarks into leptons. Reactions involving X, Y do explicitly violate both baryon number conservation and lepton number conservation since the quarks have $B = \frac{1}{3}, L_i = 0$, whereas leptons have $B = 0, L_i = 1$, where $i = e, \mu, \tau$. The baryon and lepton numbers then change, as for instance in the decay reactions

$$X \rightarrow e^- + d, \quad \Delta B = +\frac{1}{3}, \quad \Delta L_e = 1, \quad (6.105)$$

$$X \rightarrow \bar{u} + \bar{u}, \quad \Delta B = -\frac{2}{3}. \quad (6.106)$$

Second, there must be C and CP violation in the theory, as these operators change baryons into anti-baryons and leptons into anti-leptons. If the theory were C and CP symmetric, even the baryon-violating reactions (6.105) and (6.106) would be matched by equally frequently occurring reactions with opposite ΔB , so no net $B\bar{B}$ -asymmetry would result. In fact, we want baryon production to be slightly more frequent than anti-baryon production.

Third, we must require these processes to occur out of thermal equilibrium. In thermal equilibrium there is no net production of baryon number, because the X-reactions (6.105) and (6.106) go as frequently in the opposite direction. Hence the propitious moment is the phase transition when the X-bosons are freezing out of thermal equilibrium and decay. If we consult the timetable in Section 6.2, this would happen at about 10^{14} GeV: the moment for the phase transition from the GUT symmetry to its spontaneously broken remainder.

The GUT symmetry offers a good example, which we shall make use of in this section, but it is by no means obvious that GUT is the symmetry we need and that the phase transition takes place at GUT temperature. It is more likely that we have the breaking of a symmetry at a lower energy, such as supersymmetry.

Leptoquark Thermodynamics. Assuming the GUT symmetry, the scenario is therefore the following. At some energy $E_X = kT_X$ which is of the order of the rest masses of the leptoquark bosons X,

$$E_X \simeq M_X c^2, \quad (6.107)$$

all the X, Y vector bosons, the Higgs bosons, and the gluons are in thermal equilibrium with the leptons and quarks. The number density of each particle species is about the same as the photon number density, and the relations in Equation (6.100) hold.

When the age of the Universe is still young, as measured in Hubble time τ_H , compared with the mean life $\tau_X = \Gamma_X^{-1}$ of the X bosons, there are no X decays and therefore no net baryon production. The X bosons start to decay when

$$\Gamma_X \lesssim \tau_H^{-1} = H. \quad (6.108)$$

This is just like the condition in Equation (5.62) for the decoupling of neutrinos. The decay rate Γ_X is proportional to the mass M_X ,

$$\Gamma_X = \alpha M_X, \quad (6.109)$$

where α is essentially the coupling strength of the GUT interaction. It depends on the details of the GUT and the properties of the X boson.

We next take the temperature dependence of the expansion rate H from Equations (6.43) and (6.45). Replacing the Newtonian constant G by its expression in terms of the Planck mass M_P , as given in Equation (6.2), we find

$$H = \sqrt{\frac{4\pi\hbar a}{3c} g_*(T)} \frac{T^2}{M_P}. \quad (6.110)$$

Substituting this H and the expression in Equation (6.109) into Equation (6.108), the Universe is out of equilibrium when

$$AM_X \lesssim \sqrt{g_*(T)} \frac{T^2}{M_P} \quad \text{at } T = M_X, \quad (6.111)$$

where all the constants have been lumped into A . Solving for the temperature squared, we find

$$T^2 \gtrsim \frac{AM_X M_P}{\sqrt{g_*(T)}}. \quad (6.112)$$

At temperature T_X , the effective degrees of freedom g_* are approximately 100. The condition in Equation (6.112) then gives a lower limit to the X boson mass,

$$M_X \gtrsim A' \frac{M_P}{\sqrt{g_*(T)}} = A' \frac{1.2 \times 10^{19} \text{ GeV}}{\sqrt{10 \cdot 675}} \simeq A' \times 10^{18} \text{ GeV}, \quad (6.113)$$

where A' includes all constants not cited explicitly.

Thus, if the mass M_X is heavier than $A' \times 10^{18}$ GeV, the X bosons are stable at energies above M_X . Let us assume that this is the case. As the energy drops below M_X , the X and \bar{X} bosons start to decay, producing the net baryon number required. The interactions must be such that the decays really take place out of equilibrium, that is, the temperature of decoupling should be above M_X . Typically, bosons decouple from annihilation at about $M_X/20$, so it is not trivial to satisfy this requirement.

Let us now see how C and CP violation can be invoked to produce a net $\overline{B\bar{B}}$ -asymmetry in X and \bar{X} decays. The effect of the discrete operator C called *charge conjugation* on a particle state is to reverse all flavours, lepton numbers and the baryon number. The mirror reflection in three-space is called *parity transformation*, and the corresponding *parity operator* is denoted P. Obviously, every vector \mathbf{v} in a right-handed coordinate system is transformed by P into its negative in a left-handed coordinate system, CP transforms left-handed leptons into right-handed antileptons.

We can limit ourselves to the case when the only decay channels are Equations (6.105) and (6.106), and correspondingly for the \bar{X} channels. For these channels we tabulate in Table A.7 the net baryon number change ΔB and the i th branching fractions $\Gamma(X \rightarrow \text{channel } i)/\Gamma(X \rightarrow \text{all channels})$ in terms of two unknown parameters r and \bar{r} .

The baryon number produced in the decay of one pair of X, \bar{X} vector bosons weighted by the different branching fractions is then

$$\Delta B = r\Delta B_1 + (1-r)\Delta B_2 + \bar{r}\Delta B_3 + (1-\bar{r})\Delta B_4 = \bar{r} - r. \quad (6.114)$$

If C and CP symmetry are violated, r and \bar{r} are different, and we obtain the desired result $\Delta B \neq 0$. Similar arguments can be made for the production of a net lepton-anti-lepton asymmetry, but nothing is yet known about leptonic CP violation.

Suppose that the number density of X and \bar{X} bosons is N_X . We now want to generate a net baryon number density

$$N_B = \Delta B N_X \simeq \Delta B N_\gamma$$

by the time the Universe has cooled through the phase transition at T_{GUT} . After that the baryon number is absolutely conserved and further decrease in N_{B} only follows the expansion. However, the photons are also bosons, so their absolute number is not conserved and the value of η may be changing somewhat. Thus, if we want to confront the baryon production ΔB required at T_{GUT} with the present-day value of η , a more useful quantity is the baryon number per unit entropy N_{B}/S . Recall that the entropy density of photons is

$$s = 1.80g_*(T)N_\gamma \quad (6.115)$$

from Equation (6.60). At temperature T_{GUT} the effective degrees of freedom were shown in Equation (6.46) to be 106.75 (in the standard model, not counting leptoquark degrees of freedom), so the baryon number per unit entropy is

$$\frac{N_{\text{B}}}{S} = \frac{\Delta B}{1.80g_*(T_{\text{GUT}})} \simeq \frac{\Delta B}{180}. \quad (6.116)$$

Clearly this ratio scales with $g_*^{-1}(T)$. Thus, to observe a present-day value of η at about the expected value the GUT should be chosen such that it yields

$$\Delta B = \frac{g_*(T_{\text{GUT}})}{g_*(T_0)}\eta \simeq \frac{106.75}{3.36}\eta \approx 1.9 \times 10^{-8}, \quad (6.117)$$

making use of the $g_*(T)$ values in Equations (6.67) and (6.46). This is within the possibilities of various GUTs.

One may of course object that this solution of the *baryosynthesis* problem is only speculative, since it rests on the assumption that nature exhibits a suitable symmetry. At the beginning of this section, we warned that the GUT symmetry did not necessarily offer the best phase transition mechanism for baryosynthesis. The three conditions referred to could perhaps be met at some later phase transition. The reason why the GUT fails is to be found in the scenarios of cosmic inflation (Chapter 7). The baryon asymmetry produced at T_{GUT} is subsequently washed out when the Universe reheats to T_{GUT} at the end of inflation.

The search for another mechanism has turned to the electroweak phase transition at about 100 GeV. The ‘minimal standard model’ of electroweak interactions cannot generate an asymmetry but, if the correct electroweak theory could be more general. New possibilities arise if all three neutrino species oscillate and violate CP, or if one turns to the ‘minimal supersymmetric standard model’. At the expanding interface of the broken symmetry phase, the baryon–anti-baryon asymmetry could be generated via complex CP-violating reflections, transmissions and interference phenomena between fermionic excitations. Thus the existence of baryons is an indication that physics indeed has to go beyond the ‘minimal standard model’.

CPT Symmetry and Antigravity. A third discrete symmetry of importance is *time reversal* T, or symmetry under inversion of the arrow of time. This is a mirror symmetry with respect to the time axis, just as parity was a mirror symmetry with respect to the space axes. All physical laws of reversible processes are formulated in such a way that the replacement of time t by $-t$ has no observable effect. The particle reactions in

Section 5.3 occur at the same rate in both directions of the arrow. To show this, one still has to compensate for differences in phase space, i.e. the bookkeeping of energy in endothermic and exothermic reactions.

Although time reversal is not very important in itself, for instance particles do not carry a conserved quantum number related to T, it is one factor in the very important combined symmetry CPT. According to our most basic notions in theoretical physics, CPT symmetry must be absolute. It then follows from the fact that CP is not an absolute symmetry, but slightly violated, that T must be violated by an equal and opposite amount.

In a particle reaction, CPT symmetry implies that a *left-handed particle entering* the interaction region from the x -direction is equivalent to a *right-handed antiparticle leaving* the region in the x -direction. One consequence of this is that particles and antiparticles must have exactly the same mass and, if they are unstable, they must also have exactly the same mean life.

Needless to say, many ingenious experiments have been and still are carried out to test CP violation and T violation to ever higher precision. CPT symmetry will probably be tested when sufficient numbers of anti-hydrogen atoms have been produced, and their properties will be compared with those of hydrogen.

Let us assume that general relativity is CPT invariant. It certainly is so in flat space-time even if it has not been demonstrated in curved space-time. Although C transforms particles into antiparticles, invariance under CPT implies that an additional PT transformation is needed, too. This part changes the sign of each component of a four-vector and odd-rank tensor, not affecting even-rank tensors. C is ineffective because, contrary to electrodynamics, there is no charge in the gravitational field.

When CPT is applied even-rank (odd-rank) tensors which are PT-even (PT-odd) become CPT-odd (CPT-even). The change of sign of any dx^μ implies that the dt of antimatter will be reversed with respect to that of matter. Consequently the Lorentz factor $\gamma = dt/d\tau$ will be negative for antimatter while a matter particle has $\gamma = 1$ in its rest frame. Thus antiparticles are just particles travelling backwards in time. In Einsteins equation there are only even-rank tensors and scalars, so it is even under PT and CPT.

CPT invariance then implies that antimatter is attracted by antimatter in exactly the same way as matter is attracted by matter, but it is not obvious what the interaction between matter and antimatter is. In comparison with electrodynamics, the gravitational field possesses no charge, but the four-momentum for a particle of mass m

$$p^\mu = m \frac{dx^\mu}{d\tau}, \quad (6.118)$$

takes the role of a charge. Recall the expression for the four-acceleration in Equation (3.14)

$$\frac{d^2x^\mu}{d\tau^2} = -\Gamma_{\sigma\nu}^\mu \frac{dx^\sigma}{d\tau} \frac{dx^\nu}{d\tau} = 0. \quad (6.119)$$

Under (C)PT this is even because the product of the four-momenta of two particles is even, the product of the four-momenta of two antiparticles is even, the affine connection is odd, and there is a minus-sign on the right.

On the other hand, following the observation of Villata [6], if one of the momenta refers to a particle and the other one to an antiparticle, the momentum vector of the antiparticle must be (C)PT-transformed, thus acquiring a minus-sign. This then reverses the sign of the four-acceleration so that it becomes repulsive. As he points out, there are two possible interpretations for the existence of antimatter. The more conventional one is that antiparticles really exist as entities distinct from their matter counterparts, and that they travel forward in time, as all ordinary particles. The other interpretation is that antiparticles do not really exist as distinct particles, but that they are nothing else than the corresponding particles that are traveling backwards in time (BIT).

The most convincing argument in favor of the BIT interpretation is just the CPT symmetry of physical laws, since it offers a physical explanation to the need of coupling C with T (and P) for describing the behavior of antimatter: if antimatter were not traveling back in time, why should we apply the time inversion? In other words, according to the BIT interpretation, CPT is the operation that transforms events, particles and fields from one time direction to the other: the role of T is obvious, since it inverts time intervals, P is needed to get a proper Lorentz transformation (T and P alone are improper), and C provides the needed charge reversals to see time-reversed matter as antimatter.

Problems

1. Show that an expansion by a factor a leaves the black-body spectrum (6.10) unchanged, except that T decreases to T/a .
2. The flow of total energy received on Earth from the Sun is expressed by the *solar constant* $1.36 \times 10^3 \text{ J m}^{-2} \text{ s}$. Use Equation (6.41) to determine the surface temperature of the Sun,

$$b = \lambda T. \quad (6.120)$$

Using this temperature and the knowledge that the dominant color of the Sun is yellow with a wavelength of $\lambda = 0.503 \mu\text{m}$. What energy density does that flow correspond to?

3. A line in the spectrum of hydrogen has frequency $\nu = 2.5 \times 10^{15} \text{ Hz}$. If this radiation is emitted by hydrogen on the surface of a star where the temperature is 6000 K, what is the Doppler broadening [7]?
4. A spherical satellite of radius r painted black, travels around the Sun at a distance d from the center. The Sun radiates as a black-body at a temperature of 6000 K. If the Sun subtends an angle of θ radians as seen from the satellite (with $\theta \ll 1$), find an expression for the equilibrium temperature of the satellite in terms of θ . To proceed, calculate the energy absorbed by the satellite, and the energy radiated per unit time [7].
5. Use the laws of conservation of energy and momentum and the equation of relativistic kinematics [Equation (3.9)] to show that positronium cannot decay into a single photon.

6. Use the equation of relativistic kinematics [Equation (3.9)] to calculate the energy and velocity of the muon from the decay [Equation (6.51)] of a pion at rest. The neutrino can be considered massless.
7. Calculate the energy density represented by the mass of all the electrons in the Universe at the time of photon reheating when the kinetic energy of electrons is 0.2 MeV.
8. When the pions disappear below 140 MeV because of annihilation and decay, some reheating of the remaining particles occurs due to entropy conservation. Calculate the temperature-increase factor.
9. Use the equation of relativistic kinematics [Equation (3.9)] and the conservation of four-momentum to calculate the energy of the photon liberated in Equation (6.80), assuming that the ${}^4\text{He}$ nucleus is produced at rest (i.e., $v_p = v_t = v_{\text{He}} = 0$).
10. Free nucleons are favored over deuterons down to a radiation energy of 0.07 MeV. What is the ratio of photons with energies exceeding the deuteron binding energy 2.22 MeV to the number of protons at 0.07 MeV?
11. Propose a two-stage fusion process leading to the production of ${}^{12}\text{C}$.
12. Gamow's penetration factor [Equation (6.93)] gives a rough idea about the ignition temperatures in stellar interiors for each fusion reaction. Estimate these under the simplifying assumption that the burning rates during the different periods are inversely proportional to the time spans (given at the end of this chapter). Take the hydrogen burning temperature to be 10^4 K.
13. Derive a value of weak hypercharge $Y = B - L$ for the X boson from the reactions in Equations (6.72) and (6.73).

References

- [1] Coleman, T. S. and Roos, M. 2003 *Phys. Rev. D* **68**, 027702.
- [2] Peacock, J. A. 1999 *Cosmological physics*. Cambridge University Press, Cambridge.
- [3] Kolb, E. W. and Turner, M. S. 1990 *The early Universe*. Addison-Wesley, Reading, MA.
- [4] Beringer, J. *et al.* 2012 *Phys. Rev. D* **86**, Part I, 010001.
- [5] Burles, S., Nollett, K. M. and Turner, M. S. 2001 *Astrophys. J.* **552**, L1.
- [6] Villata, M. 2011 *Europhysics Letters* 94, 20001 and 2012 *Astrophys. Space Sci.* **337**, 15.
- [7] Gasiorowicz, S. 1979 *The structure of matter*. Addison-Wesley, Reading, MA.
- [8] Planck Collaboration: Ade, P. A. R. *et al.* 2014 *Astron. Astrophys.*; and Preprint arXiv:1303.5076 [astro-ph.CO]
- [9] Izotov, Y. I., Thuan T. X. and Guseva N. G. 2014, *Mon. Not. R. Astron. Soc.*; and Preprint arXiv: 1408.6953 [astro-ph.CO]

Cosmic Inflation

The concordance FLRW Big Bang model describes a homogeneous and isotropic adiabatically expanding Universe, having a beginning of space and time with very problematic initial conditions: nearly infinite temperature and density, even more homogeneous than now since inhomogeneities are unstable because of gravitation, and tend to grow with time. Where did the 10^{90} particles which make up the visible Universe come from?

Aside from such shortcomings the FLRW model has, as so far presented, been essentially a success story. We are now going to correct that optimistic picture and present a remedy: cosmic inflation.

In Section 7.1 we shall discuss problems caused by the expansion of space-time: the *horizon problem* related to its size at different epochs, the *monopole problem* associated with possible topological defects, and the *flatness problem* associated with its metric.

There are a large number of models of cosmic inflation, the simplest being the classical or *consensus* slow-roll model based on a single scalar field, which we shall study in Section 7.2. It is characterized by a de Sitter-like expansion, terminating with a huge entropy increase, in violation of the law of entropy conservation.

Most other models do not provide clear predictions regarding the nature of matter created after inflation nor the mode of exiting inflation in a vacuum that can excite the Standard Model degrees of freedom.

In Section 7.3 we discuss the scenario of the *chaotic model*, which introduces a bubble universe where we inhabit one bubble, totally unaware of other bubbles. The inflationary mechanism is the same in each bubble, but different parameter values may produce totally different universes. Since our bubble must be just right for us to exist in, this model is a version of the *Anthropic Principle*. We close this section with a discussion of the predictions of inflation.

In Section 7.4 we turn our attention to an alternative to consensus inflation, a cyclic or bouncing universe of five dimensions containing dark energy and gravity as driving forces.

7.1 Paradoxes of the Expansion

Particle Horizons. Recall the definition of the particle horizon, Equation (2.48), which in a spatially flat metric is

$$\chi_{\text{ph}} = \sigma_{\text{ph}} = c \int_{t_{\text{min}}}^{t_0} \frac{dt}{a(t)} = c \int_{a_{\text{min}}}^0 \frac{da}{a\dot{a}}. \quad (7.1)$$

This was illustrated in Figure 2.1. In expanding Friedmann models, the particle horizon is finite. Let us go back to the derivation of the time dependence of the scale factor $a(t)$ in Equations (5.39)–(5.41). At very early times, the mass density term in the Friedmann Equation (5.4) dominates over the curvature term (we have also called it the vacuum-energy term),

$$\frac{kc^2}{a^2} \ll \frac{8\pi G}{3} \rho. \quad (7.2)$$

This permits us to drop the curvature term and solve for the Hubble parameter,

$$\frac{\dot{a}}{a} = H(t) = \left(\frac{8\pi G}{3} \rho \right)^{1/2}. \quad (7.3)$$

Substituting this relation into Equation (7.1) we obtain

$$\sigma_{\text{ph}} = c \int_{a_{\text{min}}}^{a_0} \frac{da}{a^2(\dot{a}/a)} = \left(\frac{3c^2}{8\pi G} \right)^{1/2} \int_{a_{\text{min}}}^{a_0} \frac{da}{a^2 \sqrt{\rho}}. \quad (7.4)$$

In a radiation-dominated Universe, ρ scales like a^{-4} , so the integral on the right converges in the lower limit $a_{\text{min}} = 0$, and the result is that the particle horizon is finite:

$$\sigma_{\text{ph}} \propto \int_0^{a_0} \frac{da}{a^2 a^{-2}} = a_0. \quad (7.5)$$

Similarly, in a matter-dominated Universe, ρ scales like a^{-3} , so the integral also converges, now yielding $\sqrt{a_0}$. Note that an observer living at a time $t_1 < t_0$ would see a smaller particle horizon, $a_1 < a_0$, in a radiation-dominated Universe or $\sqrt{R_1} < \sqrt{R_0}$ in a matter-dominated Universe.

Suppose however, that the curvature term or a cosmological constant dominates the Friedmann equation at some epoch. Then the conditions in Equations (5.35) and (5.36) are not fulfilled; on the contrary, we have a negative net pressure

$$p < -\frac{1}{3} \rho c^2. \quad (7.6)$$

Substituting this into the law of energy conservation [Equation (5.24)] we find

$$\dot{\rho} < -2 \frac{\dot{R}}{R} \rho. \quad (7.7)$$

This can easily be integrated to give the R dependence of ρ ,

$$\rho < R^{-2}. \quad (7.8)$$

Inserting this dependence into the integral on the right-hand side of Equation (7.4) we find

$$\sigma_{\text{ph}} \propto \int_{a_{\text{min}}}^{a_0} \frac{da}{a^2 \sqrt{a^{-2}}} = \int_{a_{\text{min}}}^{a_0} \frac{da}{a}, \quad (7.9)$$

an integral which does not converge at the limit $a_{\text{min}} = 0$. Thus the particle horizon is not finite in this case. But it is still true that an observer living at a time $t_1 < t_0$ would see a particle horizon that is smaller by $\ln a_0 - \ln a_1$.

Horizon Problem. A consequence of the finite age t_0 of the Universe is that the particle horizon today is finite and larger than at any earlier time t_1 . Also, the spatial width of the past light cone has grown in proportion to the longer time perspective. Thus the spatial extent of the Universe is larger than that our past light cone encloses today; with time we will become causally connected with new regions as they move in across our horizon. This renders the question of the full size of the whole Universe meaningless—the only meaningful size being the diameter of its horizon at a given time.

In Chapter 6 we argued that thermal equilibrium could be established throughout the Universe during the radiation era because photons could traverse the whole Universe and interactions could take place in a time much shorter than a Hubble time. However, there is a snag to this argument: the conditions at any space-time point can only be influenced by events within its past light cone, and the size of the past light cone at the time of last scattering (t_{LSS}) would appear to be far too small to allow the currently observable Universe to come into thermal equilibrium.

Since the time of last scattering, the particle horizon has grown with the expansion in proportion to the $\frac{2}{3}$ -power of time (actually this power law has been valid since the beginning of matter domination at t_{eq} , but t_{LSS} and t_{eq} are nearly simultaneous). The net effect is that the particle horizon we see today covers regions which were causally disconnected at earlier times.

At the time of last scattering, the Universe was about 1090 times smaller than it is now ($z_{\text{LSS}} \approx 1065$), and the time perspective back to the Big Bang was only the fraction $t_{\text{LSS}}/t_0 \approx 2.3 \times 10^{-5}$ of our perspective. The last scattering surface is now at comoving radius 14 Gpc, but at the epoch of recombination it was at a $1 + z = 1089$ times smaller radial distance. If we assume that the Universe was radiation dominated for all the time prior to t_{LSS} , then, from Equation (5.40), $R(t) \propto \sqrt{t}$. The particle horizon at the LSS, σ_{ph} , is obtained by substituting $a(t) \propto (t_{\text{LSS}}/t)^{-1/2}$ into Equation (2.38) and integrating from zero time to t_{LSS} :

$$\sigma_{\text{ph}} \propto \int_0^{t_{\text{LSS}}} dt \left(\frac{t_{\text{LSS}}}{t} \right)^{1/2} = 2t_{\text{LSS}}. \quad (7.10)$$

It is not very critical what we call ‘zero’ time: the lower limit of the integrand has essentially no effect even if it is chosen as late as $10^{-4}t_{\text{LSS}}$.

The event horizon at the time of last scattering, σ_{eh} , represents the extent of the Universe we can observe today as light from the LSS (cf. Figures 2.1 and 7.1), since we can observe no light from before the LSS. On the other hand, the particle horizon σ_{ph}

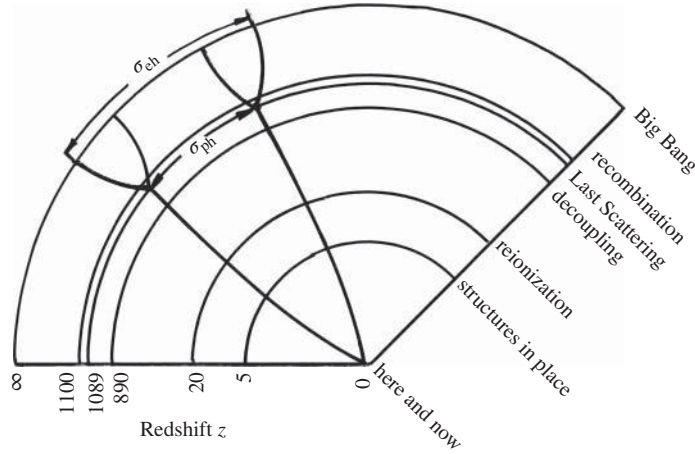


Figure 7.1 A co-moving space/conformal time diagram of the Big Bang. The observer (here and now) is at the center. The Big Bang singularity has receded to the outermost dashed circle, and the horizon scale is schematically indicated at last scattering. It corresponds to an arc of angle θ today. Reproduced from [1] by permission of J. Silk and Macmillan Magazines Ltd.

represents the extent of the LSS that could have come into causal contact from $t = 0$ to t_{LSS} . If the event horizon is larger than the particle horizon, then all the Universe we now see (in particular the relic CMB) could not have been in causal contact by t_{LSS} .

The event horizon σ_{eh} , is obtained by substituting $a(t) \propto (t_{\text{LSS}}/t)^{-2/3}$ from Equation (5.39) into Equation (2.50) and integrating over the full epoch of matter domination from t_{LSS} to $t_{\text{max}} = t_0$. Assuming flat space, $k = 0$, we have

$$\sigma_{\text{eh}} \propto \int_{t_{\text{LSS}}}^{t_0} dt \left(\frac{t_{\text{LSS}}}{t} \right)^{2/3} = 3t_{\text{LSS}} \left[\left(\frac{t_0}{t_{\text{LSS}}} \right)^{1/3} - 1 \right]. \quad (7.11)$$

Let us take $t_{\text{LSS}} = 0.35$ Myr and $t_0 = 15$ Gyr. Then the LSS particle horizon σ_{ph} is seen today as an arc on the periphery of our particle horizon, subtending an angle

$$\frac{180}{\pi} \left[\frac{\sigma_{\text{ph}}}{\sigma_{\text{eh}}} \right]_{\text{LSS}} \simeq 1.12^\circ. \quad (7.12)$$

This is illustrated in Figure 7.1, which, needless to say, is not drawn to scale. It follows that the temperature of the CMB radiation coming from any 1° arc could not have been causally connected to the temperature on a neighboring arc, so there is no reason why they should be equal. Yet the Universe is homogeneous and isotropic over the full 360° .

This problem can be avoided, as one sees from Equations (7.6)–(7.9), when the net pressure is negative, for example, when a cosmological constant dominates. In such a case, $a(t) \propto e^{\text{const.} \cdot t}$ [the case $w = -1$ in Equation (5.38)]. If a cosmological constant dominates for a finite period, say between t_1 and $t_2 < t_{\text{LSS}}$, then a term $e^{\text{const.} \cdot (t_2 - t_1)}$ enters into Equation (7.10). This term can be large, allowing a reordering of horizons to give $\sigma_{\text{ph}} > \sigma_{\text{eh}}$.

The age of the Universe at temperature 20 MeV was $t = 2$ ms and the distance scale $2ct$. The amount of matter inside that horizon was only about $10^{-5} M_\odot$, which is very

far from what we see today: matter is separated into galaxies of mass $10^{12} M_{\odot}$. The size of present superclusters is so large that their mass must have been assembled from vast regions of the Universe which were outside the particle horizon at $t = 2$ ms. But then they must have been formed quite recently, in contradiction to the age of the quasars and galaxies they contain. This paradox is the horizon problem.

The lesson of Equations (7.4)–(7.9) is that we can get rid of the horizon problem by choosing physical conditions where the net pressure is negative, either by having a large curvature term or a dominating cosmological term or some large scalar field which acts as an effective cosmological term. We turn to the latter case in Section 7.2.

GUT Phase Transition. Even more serious problems emerge as we approach very early times. At GUT time, the temperature of the cosmic background radiation was $T_{\text{GUT}} \simeq 1.2 \times 10^{28}$ K, or a factor

$$\frac{T_{\text{GUT}}}{T_0} \simeq 4.4 \times 10^{27}$$

greater than today. This is the factor by which the linear scale $a(t)$ has increased since the time t_{GUT} . If we take the present Universe to be of size $2000h^{-1}$ Mpc = 6×10^{25} m, its linear size was only 2 cm at GUT time.

Note, however, that linear size and horizon are two different things. The horizon size depends on the time perspective back to some earlier time. Thus the particle horizon today has increased since t_{GUT} by almost the square of the linear scale factor, or by

$$\frac{t_0}{t_{\text{GUT}}} = \left(\frac{g_*(T_{\text{GUT}})}{g_*(T_0)} \right)^{1/2} \left(\frac{T_{\text{GUT}}}{T_{\text{LSS}}} \right)^2 \left(\frac{T_{\text{LSS}}}{T_0} \right)^{3/2} \simeq 2.6 \times 10^{54}. \quad (7.13)$$

At GUT time the particle horizon was only 2×10^{-29} m. It follows that to arrive at the present homogeneous Universe, the homogeneity at GUT time must have extended out to a distance 5×10^{26} times greater than the distance of causal contact! Why did the GUT phase transition happen simultaneously in a vast number of causally disconnected regions? Concerning even earlier times, one may ask the same question about the Big Bang. Obviously, this paradox poses a serious problem to the standard Big Bang model.

In all regions where the GUT phase transition was completed, several important parameters—such as the coupling constants, the charge of the electron, and the masses of the vector bosons and Higgs bosons—obtained values which would characterize the present Universe. Recall that the coupling constants are functions of energy, and the same is true for particle masses. One may wonder why they obtained the same value in all causally disconnected regions.

The Higgs field had to take the same value everywhere, because this is uniquely dictated by what is its ground state. But one might expect that there would be domains where the phase transition was not completed, so that certain remnant symmetries froze in. The Higgs field could then settle to different values, causing some parameter values to be different. The physics in these domains would then be different, and so

the domains would have to be separated by *domain walls*, which are *topological defects* of space-time. Such domain walls would contain enormous amounts of energy and, in isolation, they would be indestructible. Intersecting domain walls would produce other types of topological defects such as *loops* or *cosmic strings* wiggling their way through the Universe. No evidence for topological defects has been found, perhaps fortunately for us, but they may still lurk outside our horizon.

Magnetic Monopoles. A particular kind of topological defect is a *magnetic monopole*. Ordinarily we do not expect to be able to separate the north and south poles of a bar magnet into two independent particles. As is well known, cutting a bar magnet into two produces two dipole bar magnets. Maxwell's equations account for this by treating electricity and magnetism differently: there is an electric source term containing the charge e , but there is no magnetic source term. Thus free electric charges exist, but free magnetic charges do not. Stellar bodies may have large magnetic fields, but no electric fields.

Paul A. M. Dirac (1902–1984) suggested in 1931 that the quantization of the electron charge might be the consequence of the existence of at least one free magnetic monopole with magnetic charge

$$g_M = \frac{1}{2} \frac{\hbar c n}{e} \simeq 68.5en, \quad (7.14)$$

where e is the charge of the electron and n is an unspecified integer. This would then modify Maxwell's equations, rendering them symmetric with respect to electric and magnetic source terms. Free magnetic monopoles would have drastic consequences, for instance destroying stellar magnetic fields.

Without going into detail about how frequently monopoles might arise during the GUT phase transition, we assume that there could arise one monopole per ten horizon volumes

$$N_M(t_{\text{GUT}}) = 0.1 \times (2 \times 10^{-29} \text{ m})^{-3},$$

and the linear scale has grown by a factor 4.4×10^{27} . Nothing could have destroyed them except monopole–anti-monopole annihilation, so the monopole density today should be

$$N_M(t_0) \simeq 0.1 \times (4.4 \times 0.02 \text{ m})^{-3} \simeq 150 \text{ m}^{-3}. \quad (7.15)$$

This is quite a substantial number compared with the proton density which is at most 0.17 m^{-3} . Monopoles circulating in the Galaxy would take their energy from the galactic magnetic field. Since the field survives, this sets a very low limit to the monopole flux called the *Parker bound*. Experimental searches for monopoles have not yet become sensitive enough to test the Parker bound, but they are certainly in gross conflict with the above value of N_M ; the present experimental upper limit to N_M is 25 orders of magnitude smaller than N_M .

Monopoles are expected to be superheavy,

$$m_M \gtrsim \frac{m_X}{\alpha_{\text{GUT}}} \simeq 10^{16} \text{ GeV} \simeq 2 \times 10^{-11} \text{ kg}. \quad (7.16)$$

Combining this mass with the number densities in Equations (6.76) and (7.15) the density parameter of monopoles becomes

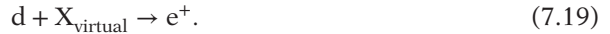
$$\Omega_M = \frac{N_M m_M}{\rho_c} \simeq 2.8 \times 10^{17}. \quad (7.17)$$

This is in flagrant conflict with the value of the dark energy density parameter $\Omega_\lambda = 1 - \Omega_m = 0.72$ to which we shall come back later. Thus, yet another paradox. Such a universe would be closed and its maximal lifetime would be only a fraction of the age of the present Universe, of the order of

$$t_{\max} = \frac{\pi}{2H_0\sqrt{\Omega_M}} \simeq 40 \text{ yr.} \quad (7.18)$$

Monopoles have other curious properties as well. Unlike the leptons and quarks, which appear to be pointlike down to the smallest distances measured (10^{-19} m), the monopoles have an internal structure. All their mass is concentrated within a core of about 10^{-30} m, with the consequence that the temperature in the core is of GUT scale or more. Outside that core there is a layer populated by the X leptoquark vector bosons, and outside that at about 10^{-17} m there is a shell of W and Z bosons.

The monopoles are so heavy that they should accumulate in the center of stars where they may collide with protons. Some protons may then occasionally penetrate in to the GUT shell and collide with a virtual leptoquark, which transforms a d quark into a lepton according to the reaction



Thus monopoles would destroy hadronic matter at a rate much higher than their natural decay rate. This would catalyze a faster disappearance of baryonic matter and yield a different timescale for the Universe.

Flatness Problem. Recall that in a spatially flat Einstein–de Sitter universe the curvature parameter k vanishes and the density parameter is $\Omega = 1$. This is obvious from Equation (5.11), where k and Ω are related by

$$\Omega - 1 = \frac{kc^2}{\dot{a}^2}.$$

The current value of the total density parameter Ω_0 is of order unity. This does not seem remarkable until one considers the extraordinary fine-tuning required: a value of Ω_0 close to, but not exactly, unity today implies that $\Omega_0(t)$ at earlier times must have been close to unity with incredible precision. During the radiation era the energy density ε_r is proportional to a^{-4} . It then follows from Equation (5.4) that

$$\dot{a}^2 \propto a^{-2}. \quad (7.20)$$

At GUT time, the linear scale was some 10^{27} times smaller than today, and since most of this change occurred during the radiation era

$$\Omega - 1 \propto a^2 \simeq 10^{-54}. \quad (7.21)$$

Thus the Universe at that time must have been flat to within 54 decimal places, a totally incredible situation. If this were not so the Universe would either have reached its maximum size within one Planck time (10^{-43} s), and thereafter collapsed into a singularity, or it would have dispersed into a vanishingly small energy density. The only natural values for Ω are therefore 0, 1 or infinity, whereas to generate a universe surviving for several Gyr without a Ω value of exactly unity requires an incredible fine-tuning. It is the task of the next sections to try to explain this.

7.2 Consensus Inflation

Let us assume that the r_p -sized universe then was pervaded by a homogeneous scalar classical field ϕ , the *inflaton* field, and that all points in this universe were causally connected. The idea with inflation is to provide a mechanism which blows up the Universe so rapidly, and to such an enormous scale, that the causal connection between its different parts is lost, yet they are similar due to their common origin. This should solve the horizon problem and dilute the monopole density to acceptable values, as well as flatten the local fluctuations to near homogeneity.

We already have tools to achieve this. In Section 5.2 on the de Sitter cosmology we saw that the solution to the FLRW Equation (5.57) for a constant expansion [Equation (5.58)] leads to an exponentially expanding universe [Equation (5.59)]. Inflationary models assume that there is a moment when the inflaton domination starts and subsequently drives the Universe into a de Sitter-like exponential expansion in which the temperature $T \simeq 0$.

Slow-roll Inflation Slow-roll inflation is a very simple idea which could be an effective representation of a variety more complicated underlying theories. It consists of one spatially homogeneous classical scalar field ϕ with a minimal kinetic term and a potential $V(\phi)$. It does not really matter what this field represents, here it is just an order parameter for a phase transition.

The total inflaton energy is of the form

$$\frac{1}{2}\dot{\phi}^2 + \frac{1}{2}(\nabla\phi)^2 + V(\phi). \quad (7.22)$$

The equation of motion for the classical field ϕ is given by the gravity model we met in Equation (5.85), now simplified with $f(R, \phi) = 1$, $\mathcal{L}_M = 0$, and with the Lagrangian

$$\mathcal{L}_\phi = \frac{1}{2}g^{\mu\nu}\partial_\mu\phi\partial_\nu\phi - V(\phi). \quad (7.23)$$

In addition we need the *Klein-Gordon* equation for the scalar field

$$\ddot{\phi} + 3H\dot{\phi} + V'(\phi) = 0. \quad (7.24)$$

Here the prime refers to derivation with respect to ϕ . The only unknown function is the potential $V(\phi)$ which contains all the important physics. Friedmann's equation becomes

$$H^2 + \frac{k}{a^2} = \frac{8\pi}{3M_p^2} \left(\frac{1}{2}\dot{\phi}^2 + \frac{1}{2}(\nabla\phi)^2 + V(\phi) \right). \quad (7.25)$$

In small-field models the field moves over a small (subPlanckian) distance. A general parametrization is the Higgs-like potential

$$V(\phi) = V_0[1 - (\phi/\mu)^p] + \dots \quad (7.26)$$

where the dots represent higher-order terms that become important near the end of inflation.

In large-field models the inflaton field starts at large field values and then evolves to a minimum at the origin $\phi = 0$. The prototypical large-field model is *chaotic inflation* where a single monomial term dominates the potential

$$V(\phi) = \lambda_p \phi^p. \quad (7.27)$$

If the field is sufficiently homogeneous, we have

$$(\nabla\phi)^2 \ll V(\phi), \quad (7.28)$$

and the $(\nabla\phi)^2$ term in Equation (7.25) then drops out.

The stress-energy tensor for a scalar field is

$$T_{\mu\nu} = (\partial_\mu\phi\partial_\nu\phi - g_{\mu\nu})\mathcal{L}_\phi, \quad (7.29)$$

and, for a homogeneous field, it takes the form of a perfect fluid with energy density

$$\rho = \frac{1}{2}\dot{\phi}^2 + V(\phi),$$

and pressure

$$p = \frac{1}{2}\dot{\phi}^2 - V(\phi).$$

In the de Sitter limit when $p \simeq -\rho$, the potential energy of the field dominates the kinetic energy, $\dot{\phi}^2 \ll V(\phi)$, and the speed of the expansion, $H = \dot{a}/a$ is large. The potential energy then acts almost as a cosmological constant $8\pi G V_0 \equiv \lambda$. Also H is almost constant and the Universe expands quasi-exponentially

$$a(t) \simeq \exp\left(\int H dt\right) \equiv e^{-N}. \quad (7.30)$$

This limit is referred to as *slow-roll*, and N is the number of e-folds that the Universe expands.

Let us rewrite the *Raychauduri Equation* (5.6) in the form

$$\frac{\ddot{a}}{a} = H^2(1 - \epsilon), \quad (7.31)$$

where the parameter ϵ specifies the Equation of State

$$\epsilon \equiv \frac{3}{2} \left(\frac{p}{\rho} + 1 \right) = \frac{4\pi G}{c^2} \left(\frac{\dot{\phi}}{H} \right)^2 = -\frac{d \ln H}{d \ln a} = H^{-1} \frac{dH}{dN}. \quad (7.32)$$

The de Sitter limit $p \simeq -\rho$ is equivalent to $\epsilon \rightarrow 0$ and accelerated expansion $a\ddot{a} > 0$ is equivalent to $\epsilon < 1$. Inflation takes place whenever $\epsilon < 1$.

If we make the further approximation that the friction term in Equation (7.24) dominates, $\ddot{\phi} \ll 3H\dot{\phi}$, the equation of motion for scalar particles is approximately

$$3H\dot{\phi} + V'(\phi) \simeq 0. \quad (7.33)$$

Slow rolling is characterised by the two parameters

$$\eta \equiv \frac{m_{\text{Planck}}^2}{16\pi} \left(\frac{V''}{V} \right) \ll 1, \quad \epsilon \equiv \frac{m_{\text{Planck}}^2}{16\pi} \left(\frac{V'}{V} \right)^2 \ll 1, \quad (7.34)$$

where $m_{\text{Planck}}^2 = c^2/G$. Single-field inflation occurs when the Universe is dominated by the inflaton field ϕ and obeys the slow-roll conditions in Equation (7.32). Inflationary models assume that there is a moment when this domination starts and subsequently drives the Universe into a de Sitter-like exponential expansion in which $T \simeq 0$. *Alan Guth* in 1981 [2] named this an *inflationary universe*.

Graceful Exit. Clearly the cosmic inflation cannot go on forever if we want to arrive at our present slowly expanding Friedmann–Lemaître universe. Thus there must be a mechanism to halt the exponential expansion, a *graceful exit*. The freedom we have to arrange this is in the choice of the potential function and its temperature-dependence, $V(\phi, T)$. Inflation ends when $\frac{1}{2}\dot{\phi}^2$ dominates over $V(\phi)$ in Friedmann's Equation (7.25) when the inflaton field arrives at the minimum $\phi = 0$ of the potential in Figure 7.2. The timescale for inflation is

$$H = \sqrt{\frac{8\pi G}{3} \langle V_0 \rangle} \propto \frac{\sqrt{\hbar c}}{M_{\text{P}}} \simeq (10^{-34} \text{ s})^{-1}. \quad (7.35)$$

One may expect that ϕ should oscillate near this minimum, but in a rapidly expanding universe, the inflaton field approaches the minimum very slowly, like a ball in a viscous medium, the viscosity $V'(\phi)$ being proportional to the speed of expansion. In the expansion the scale factor a grows so large that the third inequality follows. Equations (7.25) and (7.24) then simplify to

$$H^2 = \frac{8\pi}{3M_{\text{P}}^2} V(\phi) \quad (7.36)$$

and

$$3H\dot{\phi} = -V'(\phi). \quad (7.37)$$

There are many ways to go beyond single-field slow-roll. So far we have described the *minimally coupled* action which implies that there is no direct coupling between the inflaton and the metric. We could also entertain the possibility that the Einstein-Hilbert action is modified at high energy with $f(R)$ terms or with noncanonical terms $F(R, \phi)$ as in Equation (5.85), or that more than one field is relevant during inflation. These subjects are advanced and beyond the scope of the present monograph.

Entropy. Suppose that there is a symmetry breaking phase transition from a hot G_{GUT} -symmetric phase dominated by the scalar field ϕ to a cooler G_s -symmetric phase.

Inhomogeneities of the energy density increase as the Universe develops from an ordered, homogeneous low-entropy state towards a high-entropy chaos characterized by lower-grade heat. Thus the increase in entropy defines a *preferred direction of time*, thermal equilibrium being the state of maximum uniformity and highest entropy. The very fact that thermal equilibrium is achieved at some time tells us that the Universe must have originated in a state of low entropy.

As the Universe cools through the critical temperature T_{GUT} , bubbles of the cool phase start to appear and begin to grow. If the rate of bubble nucleation is initially small the Universe supercools in the hot phase, very much like a supercooled liquid which has a state of lowest potential energy as a solid.

The lowest curve in Figure 6.1 and the curve in Figure 7.2 illustrates the final situation when the true minimum has stabilized at ϕ_0 (denoted φ_0 in the figures), and the potential energy of this true vacuum is lower than in the original *false vacuum*:

$$V(\phi_0, T_{\text{cool}}) < V(0, T_{\text{hot}}).$$

As the potential steepens, the inflaton field begins to oscillate coherently about its vacuum state at the minimum of the potential. When the phase transition from the supercooled hot phase to the cool phase finally occurs at T_{cool} the latent heat stored as vacuum energy is liberated, *reheating* the Universe and filling it with radiation and kinetic energy of ultrarelativistic massive scalar particles with positive pressure. At the same time other GUT fields present massify in the process of spontaneous symmetry breaking, suddenly filling the Universe with particles of the reheating temperature T_{R} . Precisely how this occurs is not known, but many scenarios have been proposed.

The liberated energy heats the Universe enormously, of the order of

$$\langle V_0 \rangle \simeq (kT_{\text{R}})^4, \tag{7.38}$$

from an ambient temperature

$$T_{\text{cool}} \ll T_{\text{R}}$$

to T_{R} , which is at the T_{GUT} scale. Only at this time can one talk about a hot Big Bang.

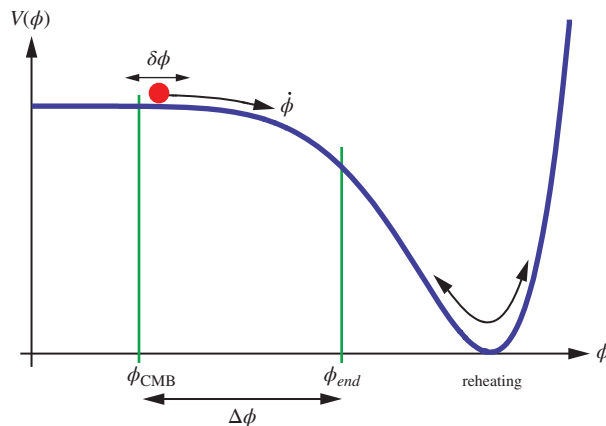


Figure 7.2 Potential energy for a real scalar field ϕ .

As long as there were no decays or annihilations of massive particles, and all other interactions conserve the total number of photons n_γ in a comoving volume is a constant. Since the entropy per photon is $s \propto T^3$ the total entropy S_H due to cosmic microwave photons of temperature T_R within our current horizon d_H is

$$S_H \approx T_R^3 d_H^3 \approx \left(\frac{T_R}{H_0} \right)^3 \approx 10^{88}. \quad (7.39)$$

the entropy per particle is suddenly increased by the very large factor

$$Z^3 = \left(\frac{T_R}{T_{\text{cool}}} \right)^3, \quad (7.40)$$

where the ratio T_R/T_{cool} is of the order of magnitude of 10^{29} . This is a very nonadiabatic process.

Using estimates of the primordial density fluctuations $\delta\rho/\rho \approx 10^{-5}$ the energy scale is of the order of $10^{-4} M_p$, so that the number of e-folds of the inflation is $N \approx 60$ or larger. There is no upper bound to N , so it could in fact be infinite, called *eternal inflation*.

At the end of inflation the Universe is a hot bubble of particles and radiation in thermal equilibrium. The energy density term in Friedmann's equations has become dominant, and the Universe henceforth follows a Friedmann–Lemaître type evolution as described in Chapters 5 and 6.

The flatness problem is now solved if the part of the Universe which became our Universe was originally homogeneous and has expanded by the de Sitter scale factor [Equation (5.59)]

$$a = e^{H\tau} \simeq 10^{29}, \quad (7.41)$$

or $H\tau \simeq 65$. Superimposed on the homogeneity of the pre-inflationary universe there were small perturbations in the field φ or in the vacuum energy. At the end of inflation these give rise to density perturbations which are the seeds of later mass structures and which can easily explain 10^{90} particles in the Universe.

It follows from Equations (7.33) and (7.36) that the duration of the inflation was

$$\tau \simeq 65 \times 10^{-34} \text{ s}. \quad (7.42)$$

Then also the horizon problem is solved, since the initial particle horizon has been blown up by a factor of 10^{29} to a size vastly larger than our present Universe. [Note that the realistic particle horizon is not infinite as one would obtain from Equation (7.9), because the lower limit of the integral is small but nonzero.] Consequently, all the large-scale structures seen today have their common origin in a microscopic part of the Universe long before the last scattering of radiation.

The development of this scenario is similar to Linde's scenario shown in Figure 7.3, except that the vertical scale here grows 'only' to 10^{29} .

When our bubble of space-time nucleated, it was separated from the surrounding supercooled hot phase by domain walls. When the phase transition finally occurred the enormous amounts of latent heat was released to these walls. The mechanism whereby this heat was transferred to particles in the bubbles was by the collision of

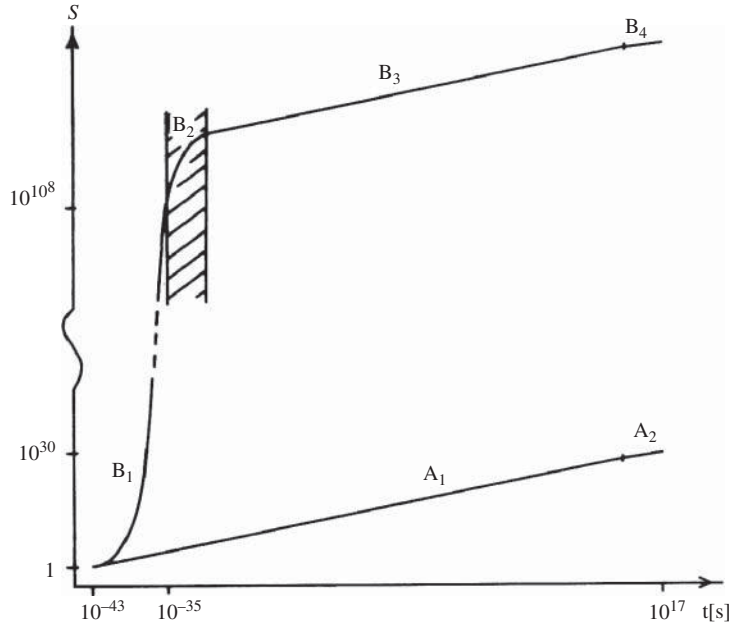


Figure 7.3 Evolution of the scale R of the Universe since Planck time in (A) Friedmann models and (B) inflationary expansion. During the epoch B_1 the Universe expands exponentially, and during B_2 the inflation ends by reheating the Universe. After the graceful exit from inflation the Universe is radiation dominated along B_3 , just as in A_1 , following a Friedmann expansion. The sections B_4 and A_2 are matter-dominated epochs.

domain walls and the coalescence of bubbles. In some models knots or topological defects then remained in the form of monopoles of the order of one per bubble. Thus the inflationary model also solves the monopole problem by blowing up the size of the region required by one monopole. There remains no inconsistency then with the present observed lack of monopoles.

Consider the contracting phase of an oscillating universe. After the time t_{\max} given by Equation (5.52) the expansion turns into contraction, and the density of matter grows. If the age of the Universe is short enough that it contains black holes which have not evaporated, they will start to coalesce at an increasing rate. Thus entropy continues to increase, so that the preferred direction of time is unchanged. Shortly before the Big Crunch, when the horizon has shrunk to linear size L_p , all matter has probably been swallowed by one enormously massive black hole.

Although Guth’s classical model of cosmic inflation may seem to solve all the problems of the hot Big Bang scenario in principle, it still fails because of difficulties with the nucleation rate. If the probability of bubble formation is large, the bubbles collide and make the Universe inhomogeneous to a much higher degree than observed. If the probability of bubble formation is small, then they never collide and there is no reheating in the Universe, so each bubble remains empty. Thus there is no graceful exit from this inflationary scenario.

Perturbations. Density perturbations in the cosmological fluid of wavelength Λ shift with the scale $a(t)$ of the expansion, or what is the same, the comoving wavelength is a constant. This is true during matter domination and radiation domination. The comoving horizon grows proportionally to the conformal time τ so that wavelengths Λ which are outside the horizon early on will later move in.

The situation was different during inflation [3]. The initial singularity is a surface of constant negative conformal time $\tau < 0$. After the Big Bang when the universe was of zero spatial size it became spatially infinite in an infinitesimal duration of time, so the Big Bang happened simultaneously everywhere at infinite speed. During inflation τ is negative, becoming less negative until $\tau = 0$ when the expansion with $a(t)$ commences. Thus perturbations of wavelengths Λ which were inside the horizon in an infinitesimal duration of time after the initial singularity will move out and disappear at superhorizon scales. The shortest wavelength perturbations are those which exited the horizon at $N = 0$ and longer ones excited earlier. Perturbations about the same size as our horizon today exited inflation at about $N = 60$.

When the perturbations enter the horizon they behave as a classical field, creating the inhomogeneities we observe in the CMB and the large-scale distribution of matter. A problem is that the separation into a homogeneous background, that depends only on time, and spatially dependent perturbations, is not unique, it depends on the choice of *gauge*.

Since the introduction of metrics in Chapter 2 we have been using the coordinates t, x, y, z or x^0, x^1, x^2, x^3 . The spacelike hyper-surfaces of constant t define the *slicing* to the four-dimensional space-time, while the timelike worldlines of constant \mathbf{x} define the *threading*. Each spacelike threading corresponds to a homogeneous universe while the slicing is orthogonal to these universes. In the homogeneous and isotropic universe we have studied so far our choice of coordinates was natural and we did not need to consider alternatives.

However, in a perturbed spacetime the definition of slicing and threading is not unique. Consider replacing the time coordinate t by a perturbed time slice $\tilde{t} = t + \delta t(t, \hat{\mathbf{x}})$. A spatially homogeneous and isotropic function is only a function of time t , as for example the energy density $\rho(t)$. On the new time-slice \tilde{t} the energy density will not be homogeneous, $\tilde{\rho} \tilde{t} \hat{\mathbf{x}} = \rho(t) (\tilde{t} \hat{\mathbf{x}})$. In general relativity we need both the matter field perturbations and the metric perturbations, so we can use the freedom of gauge transformation to trade one for the other. In the present example we can choose the hyper-surface of constant time to coincide with the hypersurface of constant energy density so that the real perturbations vanish, $\delta \tilde{\rho} = 0$.

One simple choice is to fix a gauge where the nonrelativistic limit of the full perturbed Einstein equation can be recast as a Poisson equation with a Newtonian gravitational potential, Φ . The induced metric can then be written

$$ds^2 = a^2(\tau)[(1 + 2\Phi)d\tau^2 - (1 - 2\Psi)\delta_{ik}dx^i dx^k]. \quad (7.43)$$

In the presence of Einstein gravity and when the spatial part of the energy-momentum tensor is diagonal one has $\Phi = \Psi$. Other useful gauges can be defined.

In single-field inflation we define perturbations around the homogeneous background solutions for the inflaton $\bar{\phi}$ and the metric $\bar{g}_{\mu\nu}(t)$ by

$$\phi(t, \hat{\mathbf{x}}) = \bar{\phi}(t) + \delta\phi(t, \hat{\mathbf{x}}), \quad g_{\mu\nu}(t, \hat{\mathbf{x}}) = \bar{g}_{\mu\nu}(t) + \delta g_{\mu\nu}(t, \hat{\mathbf{x}}). \quad (7.44)$$

The metric perturbations can be decomposed into scalar, vector, and tensor perturbations, of which the vectors can be ignored because they are not created by inflation and they decay with the expansion. The scalar perturbations are observed and the tensor ones will be observed as gravitational waves in the future. The ratio of tensor to scalar perturbations, denoted r , is an observable signalling gravitational waves. By the *Lyth bound* r correlates with the inflaton field moving over superPlanckian distance during inflation, provided the scalar field space is large, $> M_p$, and if also the flatness of the inflaton potential is controlled dynamically over a superPlanckian field range.

During inflation the inflationary energy is the dominant contribution to the stress-energy tensor of the universe so that the inflaton perturbations back-react on the spacetime geometry.

Primordial perturbations can give rise to nonGaussianities, but in single-field slow-roll inflation they are expected to be small. Large nonGaussianities can only arise if inflaton interactions are significant during inflation.

Initially all space-time regions of size H^{-1} would contain inhomogeneities inside their respective event horizons. At every instant during the inflationary de Sitter stage an observer would see himself surrounded by a black hole with event horizon H^{-1} (but remember that ‘black hole’ really refers to a static metric). There is an analogy between the Hawking radiation of black holes and the temperature in an expanding de Sitter space. Black holes radiate at the Hawking temperature T_H [Equation (5.83)], while an observer in de Sitter space will feel as if he is in a thermal bath of temperature $T_{\text{dS}} = H/2\pi$.

Within a time of the order of H^{-1} all inhomogeneities would have traversed the Hubble radius. Thus they would not affect the physics inside the de Sitter universe which would be getting increasingly homogeneous and flat. On the other hand, the Hubble radius is also receding exponentially, so if we want to achieve homogeneity it must not run away faster than the inhomogeneities.

The theory of perturbations during inflation is, however, beyond the scope of the present monograph. If nonGaussianities were to be observed they would be most useful to select between different inflationary scenarios.

7.3 The Chaotic Model

Consider the case of chaotic inflation with the potential

$$V(\phi) = \frac{1}{2}m_\phi^2\phi^2. \quad (7.45)$$

The time dependence of the field is then

$$\phi(t) = \phi_a - \frac{m_\phi M_P}{2\sqrt{3}\pi}t \equiv \phi_a \left(\frac{1-t}{\tau} \right), \quad (7.46)$$

where $\tau(\phi_a)$ is the characteristic timescale of the expansion. At early times when $t \ll \tau$ the scalar field remains almost constant, changing only slowly from a value $\phi_a \gg M_p$ to its ultimate value ϕ_0 . The scale factor then grows quasi-exponentially as

$$R(t) = R(t_a) \exp \left(Ht - \frac{1}{6}m_\phi^2 t^2 \right), \quad (7.47)$$

with H given by

$$H = 2\sqrt{\frac{\pi}{3}} \frac{m_\phi}{M_{\text{P}}} \phi_a. \quad (7.48)$$

At time τ , the Universe has expanded from a linear size $R(t_a)$ to

$$R(\tau) \simeq R(t_a) \exp(H\tau) = R(t_a) \exp\left(\frac{4\pi\phi_a^2}{M_{\text{P}}^2}\right). \quad (7.49)$$

For instance, a universe of linear size equal to the Planck length $R(t_a) \simeq 10^{-35}$ m has grown to

$$R(\tau) \simeq R(t_a) \exp\left(\frac{4\pi M_{\text{P}}^2}{m_\phi^2}\right). \quad (7.50)$$

For a numerical estimate we need a value for the mass m_ϕ of the inflaton. This is not known, but we can make use of the condition that the chaotic model must be able to form galaxies of the observed sizes. Then the scalar mass must be of the order of magnitude

$$m_\phi \simeq 10^{-6} M_{\text{P}}. \quad (7.51)$$

Inserting this estimate into Equation (7.43) we obtain the completely unfathomable scale

$$R(\tau) \simeq 10^{-35+\exp(4\pi\times 10^{12})} \text{ m} \simeq 10^{5.5\times 10^{12}} \text{ m}. \quad (7.52)$$

It is clear that all the problems of the standard Big Bang model discussed in Section 7.1 then disappear. The homogeneity, flatness and isotropy of the Universe turn out to be consequences of the inflaton field having been large enough in a region of size M_{P}^{-1} at time t_{p} . The inflation started in different causally connected regions of space-time ‘simultaneously’ to within 10^{-43} s, and it ended at about 10^{-35} s. Our part of that region was extremely small. Since the curvature term in Friedmann’s equations decreased exponentially, the end result is exactly as if k had been zero to start with. A picture of this scenario is shown in Figure 7.3.

Quantum Fluctuations. At Planck time the universe cannot have been completely homogeneous and isotropic because of quantum fluctuations. Across the horizon of size M_{P}^{-1} the field may have varied by an amount

$$\Delta\phi_a \simeq M_{\text{P}}. \quad (7.53)$$

But in quantum mechanics we noted that at Planck time the field ϕ was indefinite by M_{P} , at least, so that there were deviations from a pure de Sitter universe. Even if this universe was empty, quantum field theory tells us that empty space is filled with zero-point quantum fluctuations of all kinds of physical fields, here fluctuations from the classical de Sitter inflaton field.

The vacuum fluctuation spectrum of the slowly rolling scalar field during the inflationary expansion turns out to be quite unlike the usual spectrum of thermal

fluctuations. This can be seen if one transforms the de Sitter metric (5.60) into the metric of a Euclidean four sphere [Equation (2.28)]. Bose fields (like the inflaton) obeying a massless Klein–Gordon equation turn out to oscillate harmonically on this sphere with period $2\pi/H$, which is equivalent to considering quantum statistics at a temperature $T_{\text{dS}} = H/2\pi$. However, the temperature in de Sitter space is highly unusual in that the fluctuations on the four sphere are periodic in all four dimensions [4, 5].

The fate of a bubble of space-time clearly depends on the starting value of ϕ . Only when it is large enough will inflationary expansion commence. If ϕ is very much larger than M_{p} , Equation (7.41) shows that the rate of expansion is faster than the timescale τ ,

$$H \gg 2\sqrt{\frac{\pi}{3}}m_{\phi} \simeq \frac{2}{\tau}. \tag{7.54}$$

Although the wavelengths of all quantum fields then grow exponentially, the change $\Delta\phi$ in the value of the inflaton field itself may be small. In fact, when the physical wavelengths have reached the size of the Hubble radius H^{-1} , all changes in ϕ are impeded by the friction $3H\dot{\phi}$ in Equation (7.37), and fluctuations of size $\delta\phi$ freeze to an average nonvanishing amplitude of

$$|\delta\phi(x)| \simeq \frac{H}{2\pi}. \tag{7.55}$$

Consequently, the vacuum no longer appears empty and devoid of properties.

Fluctuations of a length scale exceeding H^{-1} are outside the present causal horizon so they no longer communicate, crests and troughs in each oscillation mode remain frozen. But at the end of inflation, the expansion during radiation and matter domination starts to return these frozen fluctuations inside the horizon. With time they become the seeds of perturbations we now should observe in the CMB and in the density distribution of matter.

The quantum fluctuations remaining in the inflaton field will cause the energy to be dumped into entropy at slightly fluctuating times. Thus, the Universe will also contain entropy fluctuations as seeds of later density perturbations.

Note that the quantum fluctuations amplified during inflation also lead to self-reproduction and therefore do not allow inflation to end once it has started. Inflation continues forever leading to a metaphysical (nonverifiable) concept of eternal universe and multiverse. This damages the predictive power of the theory because in this case ‘anything can happen and will happen an infinite number of times’.

Linde’s Bubble Universe. Since our part of the pre-inflationary universe was so small, it may be considered as just one bubble in a foam of bubbles having different fates. In Linde’s chaotic model each original bubble has grown in one e-folding time $\tau = H^{-1}$ to a size comprising e^3 mini-universes, each of diameter H^{-1} . In half of these mini-universes, on average, the value of ϕ may be large enough for inflation to continue, and in one-half it may be too small. In the next e-folding time the same pattern is repeated. Linde has shown that in those parts of space-time where ϕ grows

continuously the volume of space grows by a factor

$$e^{(3-\ln 2)Ht}, \quad (7.56)$$

whereas in the parts of space-time where ϕ does not decrease the volume grows by the factor

$$\frac{1}{2}e^{3Ht}. \quad (7.57)$$

Since the Hubble parameter is proportional to ϕ , most of the physical volume must come from bubbles in which ϕ is maximal:

$$\phi \simeq M_{\text{P}}^2/m_{\phi}. \quad (7.58)$$

But there must also be an exponential number of bubbles in which ϕ is smaller. Those bubbles are the possible progenitors of universes of our kind. In them, ϕ attains finally the value corresponding to the true minimum $V(\phi_0)$, and a Friedmann–Lemaître-type evolution takes over. Elsewhere the inflatoric growth continues forever. Thus we happen to live in a universe which is a minuscule part of a steady-state eternally inflating meta-universe which has no end, and therefore it also has no beginning. There is simply no need to turn inflation on in the first place, and the singularity at time zero has dropped out from the theory.

During inflation, each bubble is generating new space-time to expand into, as required by general relativity, rather than expanding into pre-existing space-time. In these de Sitter space-times the bubble wall appears to an observer as a surrounding black hole. Two such expanding bubbles are causally disconnected, so they can neither collide nor coalesce. Thus the mechanism of vacuum-energy release and transfer of heat to the particles created in the phase transition is not by bubble collisions as in the classical model. Instead, the rapid oscillations of the inflaton field ϕ decay away by particle production as the Universe settles in the true minimum. The potential energy then thermalizes and the Universe reheats to some temperature of the order of T_{GUT} .

In this reheating, any baryon–anti-baryon asymmetry produced during the GUT phase transition mechanism is washed out, that is why some other phase transition must be sought to explain the baryon–anti-baryon asymmetry. Thus the existence of baryons is an indication that particle physics indeed has to go beyond the ‘minimal standard model’.

7.4 Predictions

One consequence of the repulsive scalar field is that any two particles appear to repel each other. This is the Hubble expansion, which is a consequence of inflation. In non-inflationary theories the Hubble expansion is merely taken for granted.

Inflationary models predict that the density of the Universe should today be very nearly critical,

$$\Omega_0 = 1. \quad (7.59)$$

This prediction is verified to within 0.3 %. Consequently, we should not only observe that there is too little luminous matter to explain the dynamical behavior of the

Universe, we also have a precise theoretical specification for how much matter there should be. This links dark matter to inflation.

We have already noted that the scalar inflaton field produced a spectrum of frozen density and radiation perturbations beyond the horizon, which moved into sight when the expansion of the Universe decelerated. In the post-inflationary epoch when the Friedmann expansion takes over we can distinguish between two types of perturbations, *adiabatic fluctuations*, also called *curvature perturbations*, and *isocurvature fluctuations*, also called *isothermal* perturbations. In the first case, the perturbations in the local number density, $\delta_m \equiv \delta\rho_m/\rho_m$, of each species of matter—baryons, leptons, neutrinos, dark matter—is the same. In particular, these perturbations are coupled to those of radiation, $\delta_r \equiv \delta\rho_r/\rho_r$, so that $4\delta_m = 3\delta_r$ [from Equation (6.39)]. By the principle of covariance, perturbations in the energy-momentum tensor imply simultaneous perturbations in energy density and pressure, and by the equivalence principle, variations in the energy-momentum tensor are equivalent to variations in the curvature. Curvature perturbations can have been produced early as irregularities in the metric, and they can then have been blown up by inflation far beyond the Hubble radius. Thus adiabatic perturbations are a natural consequence of cosmic inflation. In contrast, inflation does not predict any isocurvature perturbations.

Let us write the power spectrum of density perturbations in the form

$$P(k) \propto k^{n_s}, \quad (7.60)$$

where n_s is the *scalar spectral index*. Inflationary models predict that the primordial fluctuations have an equal amplitude on all scales, an almost scale-invariant power spectrum as the matter fluctuations cross the Hubble radius, and are Gaussian. This is the Harrison–Zel’dovich spectrum for which $n_s = 1$ ($n_s = 0$ would correspond to white noise).

A further prediction of inflationary models is that tensor fluctuations in the space-time metric, satisfying a massless Klein–Gordon equation, have a nearly scale-invariant spectrum of the form in Equation (7.59) with *tensor spectral index* n_t , just like scalar density perturbations, but independently of them. The ratio $r = n_t/n_s$ is now eagerly measured by several teams.

The above predictions are generic for a majority of inflation models which differ in details. Inflation as such cannot be either proved or disproved, but specific theories can be and will be ruled out by these observations.

7.5 A Cyclic Universe

As we have seen, ‘consensus’ inflation by a single inflaton field solves the problems described in Section 7.1. But in the minds of some people it does so at a very high price. It does not explain the beginning of space and time, it does not predict the future of the Universe, or it sweeps these fundamental questions under the carpet of the *anthropic principle*. It invokes several unproven ingredients, such as a scalar field and a scalar potential, suitably chosen for the field to slow-roll down the potential while its kinetic energy is negligible, and such that it comes to a graceful exit where

ordinary matter and radiation are created by oscillations in the potential well, or by entropy generation during a second slow-roll phase of an equally arbitrary dark energy field. Clearly, any viable alternative to single-field inflation must also be able to solve the problems in Section 7.1, and it should not contain more arbitrary elements than does single-field inflation—multiple scalar fields have more freedom but also more arbitrariness.

In the radiation-dominated Universe, the source of energy of photons and other particles is a phase transition or a particle decay or an annihilation reaction, many of these sources producing monoenergetic particles. Thus the energy spectrum produced at the source is very nonuniform and nonrandom, containing sharp spectral lines ‘ordered’ by the type of reaction. Such a spectrum corresponds to low entropy. Subsequent scattering collisions will redistribute the energy more randomly and ultimately degrade it to high-entropy heat. Thermal equilibrium is the state of maximum uniformity and highest entropy. The very fact that thermal equilibrium is achieved at some time tells us that the Universe must have originated in a state of low entropy.

In the transition from radiation domination to matter domination no entropy is lost. We have seen the crucial effect of photon reheating due to entropy conservation in the decoupling of the electrons. As the Universe expands and the wavelengths of the CMB photons grow, the available energy is continuously converted into lower-grade heat, thus increasing entropy. This thermodynamic argument defines a *preferred direction of time*.

When the cooling matter starts to cluster and contract under gravity, a new phase starts. We have seen that the Friedmann–Lemaître equations dictate instability, the lumpiness of matter increases with the Universe developing from an ordered, homogeneous state towards chaos. It may seem that contracting gas clouds represent higher uniformity than matter clustered into stars and galaxies. If the only form of energy were thermal, this would indeed be so. It would then be highly improbable to find a gas cloud squeezed into a small volume if nothing hinders it from filling a large volume. However, the attractive nature of gravity seemingly reverses this logic: the cloud gains enormously in entropy by contracting. Thus the preferred direction of time as defined by the direction of increasing entropy is unchanged during the matter-dominated era.

The same trend continues in the evolution of stars. Young suns burn their fuel through a chain of fusion reactions in which energetic photons are liberated and heavier nuclei are produced. As the photons diffuse in the stellar matter, they ultimately convert their energy into a large number of low-energy photons and heat, thereby increasing entropy. Old suns may be extended low-density, high-entropy red giants or white dwarfs, without enveloping matter which loses mass by various processes. In the process of supernova explosion entropy grows enormously.

Consider the contracting phase of an oscillating universe. After the time t_{\max} given by Equation (5.52) the expansion turns into contraction, and the density of matter grows. If the age of the Universe is short enough that it contains black holes which have not evaporated, they will start to coalesce at an increasing rate. Thus entropy continues to increase, so that the preferred direction of time is unchanged. Shortly before the Big Crunch, when the horizon has shrunk to linear size L_p , all matter has probably been swallowed by one enormously massive black hole.

Early attempts to build models with cyclically reoccurring expansion and contraction were plagued by the problem, that the entropy density would rise from cycle to cycle. The length of cycles must then increase steadily. But, in retrospect, there must then have been a first cycle a finite time ago, thus a beginning of time: precisely what the cyclic model was conceived to avoid.

A cyclic model which solves the entropy problem and which appears as successful as the ‘consensus’ inflationary model has been proposed by Steinhardt and Turok [6].

The model is described qualitatively in Figure 7.4, which depicts a potential $V(\phi)$, function of a scalar field ϕ . Unlike the inflaton field, ϕ does not cause an inflationary expansion of space-time. Following the arguments of Steinhardt and Turok the universe is assumed to be a five-dimensional bulk, where our physical world with all its observable particles is located on a four-dimensional *brane* separated from another brane by a microscopic gap (or rather “Planckoscopic”). Matter on the other brane can only interact with our world gravitationally, but not through strong or electromagnetic interactions. Such matter is therefore dark, we cannot detect it in our laboratories.

The potential $V(\phi)$ describes a force between the two branes, and ϕ is a moduli field that describes the interbrane separation. When the two branes are far apart this force is very weak, the branes are stretched to the point where they are flat and parallel, and we experience this as the present slow cosmic acceleration (phase 1 in Figure 7.4), the situation where we are now. During this phase the potential energy of the scalar field dominates over the kinetic energy of the scalar field so that dark energy drives the expansion.

While ϕ decreases toward phase 2, the potential slow-rolls down the weakly sloping positive potential and the branes draw together. At point 3 in Figure 7.4 the potential

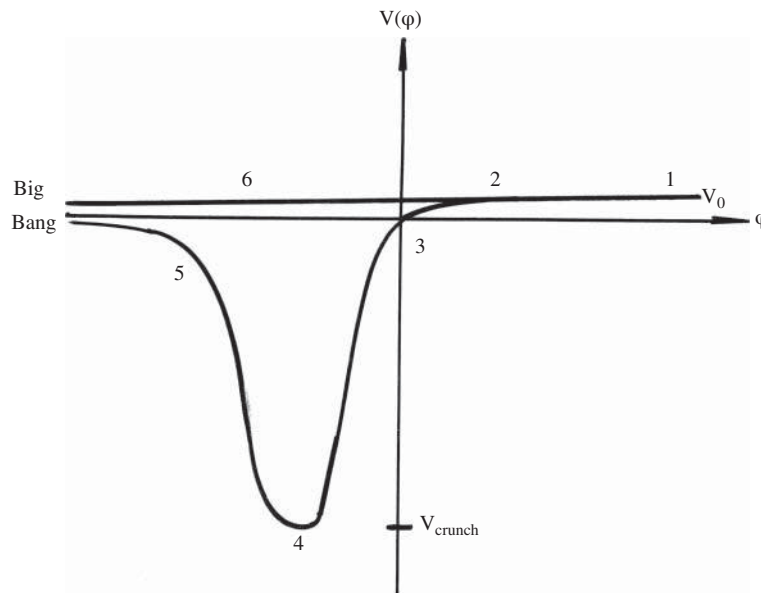


Figure 7.4 Schematic view of the potential $V(\phi)$ as a function of the field ϕ for cyclic models in a five-dimensional universe. The numbered sequence of phases is described in the text.

has decreased to the point where $V(\phi) = \frac{1}{2}\phi^2$ where the field crosses zero potential, the total energy density is zero and the Hubble expansion halts. As the branes draw closer together, the potential energy decreases from positive to negative values and the universe switches from expansion to a rapid contraction with equation of state $w \gg 1$, towards a minimum at point 4 in Figure 7.4. The spatial contraction occurs in the extra dimension ϕ , rather than in our three dimensions.

As the scalar field rolls down into this minimum (phase 4), it picks up speed on the way and causes a Hubble blueshift of the scalar field kinetic energy. With this speed the field just continues beyond the minimum, climbs out of it on the other side and continues towards the bounce at negative infinity in a finite time. The potential disappears exponentially. In the approach towards $\phi = -\infty$, the scale factor $a(t)$ goes to zero (phase 5), but the coupling of the scalar field to radiation and matter conspires in such a way as to keep the temperature and the energy density finite. Space does not disappear completely since this scenario takes place in a five-dimensional space-time, and $a(t)$ has to be interpreted as the effective scale factor on the brane.

Each cycle ends and begins with a crunch turning into a bang at the field value $\phi = -\infty$. The bang is a turn-around or bounce from the pre-existing contracting phase with a decreasing field. At the bounce, the branes collide and this is the big bang of a new cycle. After the branes bounce apart there immediately follows an expanding phase and an acceleration of the field towards positive values (phase 6 in Figure 7.4), when matter and radiation are created at large but finite temperature from the kinetic energy of the field.

This phase is quite similar to the corresponding post-inflationary epoch in the conventional inflationary scenario, and therefore the predictions are the same. But, unlike the conventional model, here a subdominant dark energy field is required. During radiation and matter domination the scalar dark energy field is effectively frozen in place by the Hubble redshift of its kinetic energy and the branes slow down to essentially a halt. The potential energy of this field starts to dominate only when radiation and matter energy densities have been sufficiently diluted by the expansion; then a slow cosmic acceleration commences (phase 1).

The heat from the collision dominates the universe for a few billion years, but eventually, on the way from phase 6 to phase 1, it is diluted enough that the positive inter-brane potential energy density dominates. This acts as a source of dark energy that causes the expansion of the branes to accelerate. Although the acceleration due to dark energy is very slow, causing the Universe to double in size every 15 billion years or so, compared with the enormous expansion in Equation (7.45) during 10^{-35} s, this is enough to empty the Universe of its matter and radiation. The matter, radiation, and large scale structure dilute away over the next trillion years or so, and the branes become nearly perfect vacua at the end of each cycle (at phase 1), preparing the way for a new cycle of identical duration.

The homogeneity and flatness of the Universe and the density perturbations are established during long periods of ultra-slow accelerated expansion, and the conditions are set up during the negative time prior to a bang. In contrast, standard inflationary theories have very little time to set up large-scale conditions, only about 10^{-35} s until the inflationary fluctuations have been amplified and frozen in, at a length

scale of 10^{-25} cm. In the cyclic Universe, the fluctuations can be generated a fraction of a second before the bang when their length scale is thousands of kilometres.

The entropy created in one cycle is expanded and exponentially diluted to near zero density after the dark energy dominated phase, but the entropy does not increase again in the contracting phase. The reason is that the branes themselves do not contract, only the extra dimension contracts. The total entropy on the branes and the number of black holes increase from cycle to cycle, and increase per comoving volume on our brane as well. In the vicinity of the black holes, there is no cycling due to their strong gravitational field. Black holes formed during one cycle will therefore survive to the next cycle, acting as defects in an otherwise nearly uniform universe.

The kinetic energy and the physical entropy density during each cycle are fed not only by the interbrane force, but also by gravity which supplies extra energy during the contraction phase. This kinetic energy of the branes is converted partially into matter and radiation and blue-shifted by the gravity.

All of this is very speculative, but so is consensus inflation, too. Fortunately, the different models make different testable predictions, notably for gravitational radiation. A certain discovery of gravitational radiation testified by tensor perturbations in the CMB would support consensus inflation.

Problems

1. Derive Equation (7.37).
2. Derive $\phi(t)$ for a potential $V(\phi) = \frac{1}{4}\lambda\phi^4$.
3. Suppose that the scalar field averaged over the Hubble radius H^{-1} fluctuates by an amount ψ . The field gradient in this fluctuation is $\nabla\psi = H\psi$ and the gradient energy density is $H^2\psi^2$. What is the energy of this fluctuation integrated over the Hubble volume? Use the timescale H^{-1} for the fluctuation to change across the volume, and the uncertainty principle to derive the minimum value of the energy. This is the amount by which the fluctuation has stretched in one expansion time [7].
4. Material observed now at redshift $z = 1$ is at present distance H_0^{-1} . The recession velocity of an object at coordinate distance x is $\dot{R}x$. Show that the recession velocity at the end of inflation is

$$\dot{R}x = \frac{H_0 R_0 x z_r}{\sqrt{z_{\text{eq}}}}, \tag{7.61}$$

where z_r is the redshift at the end of inflation. Compute this velocity. The density contrast has grown by the factor z_r^2/z_{eq} . What value did it have at the end of inflation since it is now $\delta \approx 10^{-4}$ at the Hubble radius [7]?

5. Show that in an exponentially expanding universe ($q = -1$) the Hubble sphere is stationary. Show that it constitutes an event horizon in the sense that events beyond it will never be observable. Show that in this universe there is no particle horizon [8].

6. Show that the number of e-foldings of inflation in the $V(\phi) = -\lambda\phi^4$ model is of order

$$N \approx \frac{H^2}{\lambda\phi_i^2}$$

from the time at which the field has the value ϕ_i to the end of inflation ($\phi \ll \phi_i$). Hence show that the density perturbations in this model are of order

$$\left(\frac{\delta\rho}{\rho}\right)_H \approx \sqrt{\lambda N^3}. \quad (7.62)$$

Deduce that $\lambda < 10^{-14}$ is required if the fluctuations are to be compatible with the CMB. This of course amounts to the fine-tuning that inflation is supposed to avoid [8].

References

- [1] Kaiser, N. and Silk, J. 1986 *Nature* **324**, 529.
- [2] Guth, A. H. 1981 *Phys. Rev. D* **23**, 347.
- [3] Kenney, W. H. 2009 *TASI Lectures on Inflation, Boulder* and preprint arXiv: 0902.1529 astro-ph.CO.
- [4] Linde, A. D. 1990 *Particle physics and inflationary cosmology*. Harwood Academic Publishers, London.
- [5] Linde, A. D. 2002 Inflationary cosmology and creation of matter in the Universe. In *Modern cosmology* (ed. S. Bonometto, V. Gorini and U. Moschella). Institute of Physics Publishing, Bristol.
- [6] Steinhardt, P. J. and Turok, N. 2002 *Science* **296**, 1496; 2005 *New Astron. Rev.* **49**, 43.
- [7] Peebles, P. J. E. 1993 *Principles of physical cosmology*. Princeton University Press, Princeton, NJ.
- [8] Raine, D. J. and Thomas, E. G. 2001 *An introduction to the science of cosmology*. Institute of Physics Publishing, Bristol.

8

Cosmic Microwave Background

In this chapter we shall meet several important observational discoveries. The cosmic microwave background (CMB), which is a consequence of the hot Big Bang and the following radiation-dominated epoch, was discovered in 1964. We discuss this discovery in Section 8.1.

The hot Big Bang also predicts that the CMB radiation should have a blackbody spectrum. Inflation predicts that the mean temperature of the CMB should exhibit minute perturbations across the sky. These predictions were verified by a 1990 satellite experiment, the *Cosmic Background Explorer (COBE)*. The COBE observations verified the existence of a nearly scale-invariant spectrum of primordial fluctuations on angular scales larger than 7° . Many experiments have since then extended the range, in particular the Wilkinson Microwave Anisotropy Probe *WMAP*. Most recently the third-generation space mission *Planck* has made observations of temperature and polarization anisotropies with much improved precision and statistics at different angular scales. These measurements will be discussed in the rest of this chapter.

In Section 8.2 we shall discuss the method of analyzing the temperature perturbations which are expected to be associated with the even smaller polarization variations, due to Thomson scattering at the LSS. These were first observed by the ground-based Degree Angular Scale Interferometer (DASI) in late 2002 and by in early 2003. We discuss this in Section 8.3.

The CMB contains a wealth of information about the dynamical parameters of the Universe and on specific features of the theoretical models: general relativity, the standard FLRW cosmology versus other cosmologies, all versions of inflation and its alternatives, dark matter, dark energy and so on. In Section 8.4 we establish the parameters, how they are related to each other, what observational values they have and what information they give about possible cosmological models.

8.1 The CMB Temperature

Predictions. In 1948, *Georg Gamow* (1904–1968), *Ralph Alpher* and *Robert Herman* calculated the temperature at that time of the primordial blackbody radiation which started free streaming at the LSS. They found that the CMB should still exist today, but that it would have cooled in the process of expansion to the very low temperature of $T_0 \approx 5$ K. This corresponds to a photon wavelength of

$$\lambda = \frac{hc}{kT_0} = 2.9 \times 10^{-3} \text{ m.} \quad (8.1)$$

This is in the microwave range of radio waves (see Table A.3). (The term ‘microwave’ is actually a misnomer, since it does not refer to micrometer wavelengths, but rather to centimeters.)

We can now redo their calculation, using some hindsight. Let us first recall from Equations (5.39) and (5.40) that the expansion rate changed at the moment when radiation and matter contributed equally to the energy density. For our calculation we need to know this equality time, t_{eq} , and the temperature T_{eq} . The radiation energy density is given by Equation (6.41):

$$\epsilon_r = (g_\gamma + 3g_\nu) \frac{1}{2} a T_{\text{eq}}^4. \quad (8.2)$$

The energy density of matter at time T_{eq} is given by Equation (6.24), except that the electron (e^- and e^+) energy needs to be averaged over the spectrum [Equation (6.28)]. We could in principle solve for T_{eq} by equating the radiation and matter densities,

$$\epsilon_r(T_{\text{eq}}) = \rho_m(T_{\text{eq}}). \quad (8.3)$$

We shall defer solving this to Section 8.4.

The transition epoch happens to be close to the recombination time [z_{rec} in redshift, see Equation (6.70)] and the LSS [z_{LSS} in redshift, see Equation (6.71)]. With their values for t_{eq} , T_{eq} and t_0 and Equation (5.39), Gamow, Alpher and Herman obtained a prediction for the present temperature of the CMB:

$$T_0 = T_{\text{eq}} \left(\frac{t_{\text{eq}}}{t_0} \right)^{2/3} = 2.45 \text{ K.} \quad (8.4)$$

This is very close to the present-day observed value, as we shall see.

Discovery. Nobody paid much attention to the prediction of Gamow *et al.*, because the Big Bang theory was generally considered wildly speculative, and detection of the predicted radiation was far beyond the technical capabilities existing at that time. In particular, their prediction was not known to *Arno Penzias* and *Robert Wilson* who, in 1964, were testing a sensitive antenna intended for satellite communication. They wanted to calibrate it in an environment free of all radiation, so they chose a wavelength of $\lambda = 0.0735$ m in the relatively quiet window between the known emission from the Galaxy at longer wavelengths and the emission at shorter wavelengths from the Earth’s atmosphere. They also directed the antenna high above the galactic plane, where scattered radiation from the Galaxy would be minimal.

To their consternation and annoyance they found a constant low level of background noise in every direction. This radiation did not seem to originate from distant galaxies, because in that case they would have seen an intensity peak in the direction of the nearby M31 galaxy in Andromeda. It could also not have originated in Earth's atmosphere, because such an effect would have varied with the altitude above the horizon as a function of the thickness of the atmosphere.

Thus Penzias and Wilson suspected technical problems with the antenna (in which a couple of pigeons turned out to be roosting) or with the electronics. All searches failing, they finally concluded, correctly, that the Universe was uniformly filled with an 'excess' radiation corresponding to a blackbody temperature of 3.5 K, and that this radiation was isotropic and unpolarized within their measurement precision.

At Princeton University, a group of physicists led by *Robert Dicke* (1916–1997) had at that time independently arrived at the conclusion of Gamow and collaborators, and they were preparing to measure the CMB radiation when they heard of the remarkable 3.5 K 'excess' radiation. The results of Penzias and Wilson's measurements were immediately understood and they were subsequently published (in 1965) jointly with an article by Dicke and collaborators which explained the cosmological implications. The full story is told by Peebles [1], who was a student of Dicke at that time. Penzias and Wilson (but not Gamow or Dicke) were subsequently awarded the Nobel prize in 1978 for this discovery.

This evidence for the 15-Gyr-old echo of the Big Bang counts as the most important discovery in cosmology since Hubble's law. In contrast to all radiation from astronomical bodies, which is generally hotter, and which has been emitted much later, the CMB has existed since the era of radiation domination. It is hard to understand how the CMB could have arisen without the cosmic material having once been highly compressed and exceedingly hot. There is no known mechanism at any time after decoupling that could have produced a blackbody spectrum in the microwave range, because the Universe is transparent to radio waves.

Spectrum. In principle, one intensity measurement at an arbitrary wavelength of the blackbody spectrum (6.10) is sufficient to determine its temperature, T , because this is the only free parameter. On the other hand, one needs measurements at different wavelengths to establish that the spectrum is indeed blackbody.

It is easy to see that a spectrum which was blackbody at time t with temperature T will still be blackbody at time t' when the temperature has scaled to

$$T' = T \frac{a(t)}{a(t')}. \tag{8.5}$$

This is so because, in the absence of creation or annihilation processes, the number of photons, $n_\gamma a^3(t)$, is conserved. Thus the number density $dn_\gamma(\nu)$ in the frequency interval $(\nu, \nu + d\nu)$ at time t transforms into the number density at time t' ,

$$dn'_\gamma(\nu') = \left[\frac{a(t)}{a(t')} \right]^3 dn_\gamma(\nu). \tag{8.6}$$

Making use of Equations (6.10) and (8.5), the distribution at time t' becomes

$$dn'_\gamma(\nu') = \frac{8\pi}{3c^3} \frac{d\nu'^3}{e^{h\nu'/kT'} - 1}, \quad (8.7)$$

which is precisely the blackbody spectrum at temperature T' .

Although several accurate experiments since Penzias and Wilson have confirmed the temperature to be near 3 K by measurements at different wavelengths, the conclusive demonstration that the spectrum is indeed also blackbody in the region masked by radiation from the Earth's atmosphere was first made by a dedicated instrument, the Far Infrared Absolute Spectrophotometer (FIRAS) aboard the COBE satellite launched in 1989 [2]. The present value of the photon temperature, T_0 , is from measurements in 2009,

$$T_0 = 2.7255 \pm 0.0006 \text{ K}. \quad (8.8)$$

The spectrum reported by the COBE team in 1993 [3] matches exactly the theoretical prediction for blackbody radiation remaining from the hot Big Bang. The measurement errors on each of the 34 wavelength positions have not been added to the theoretical blackbody curve because they could not be distinguished. It is worth noting that such a pure blackbody spectrum had never been observed in laboratory experiments. All theories that attempt to explain the origin of large-scale structure seen in the Universe today must conform to the constraints imposed by these COBE measurements.

The vertical scale in Figure 8.1 gives the *intensity* $I(1/\lambda)$ of the radiation, that is, the power per unit inverse wavelength interval arriving per unit area at the observer from

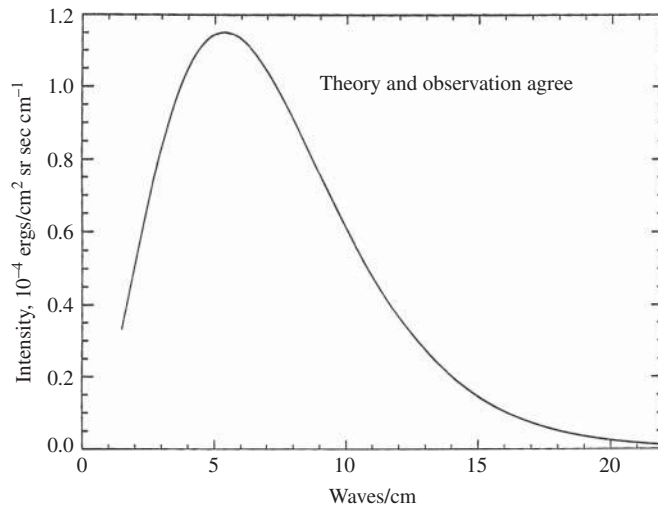


Figure 8.1 The theoretical blackbody spectrum from Equation (6.13). From D. J. Fixsen, E. S. Cheng, J. M. Gales, J. C. Mather, R. A. Shafer, and E. L. Wright, *The Cosmic Microwave Background Spectrum from the Full COBE* FIRAS Data Set*. *Astrophys. J.*, **473**, 576, published 20 December ©1996. AAS. Reproduced with permission.

one steradian of sky. In SI units this is $10^{-9} \text{ J m}^{-1} \text{ sr}^{-1} \text{ s}^{-1}$. This quantity is equivalent to the intensity per unit frequency interval, $I(\nu)$. One can transform from $d\lambda$ to $d\nu$ by noting that $I(\nu) d\nu = I(\lambda) d\lambda$, from which

$$I(\lambda) = \frac{\nu^2}{c} I(\nu). \quad (8.9)$$

The relation between energy density ϵ_r and total intensity, integrated over the spectrum, is

$$\epsilon_r = \frac{4\pi}{c} \int I(\nu) d\nu. \quad (8.10)$$

Energy and Entropy Density. Given the precise value of T_0 in Equation (8.8), one can determine several important quantities. From Equation (6.41) one can calculate the present energy density of radiation

$$\epsilon_{r,0} = \frac{1}{2} g_{*S} a_S T_0^4 = 2.606 \times 10^5 \text{ eV m}^{-3}. \quad (8.11)$$

The corresponding density parameter then has the value

$$\Omega_r = \frac{\epsilon_{r,0}}{\rho_c} = 2.473 \times 10^{-5} \text{ h}^{-2}, \quad (8.12)$$

using the value of ρ_c from Equation (1.31) and the compromise value $H_0 = 0.71$. Obviously, the radiation energy is very small today and far from the value $\Omega_0 = 1$ required to close the Universe.

The present value of the entropy density is

$$s = \frac{4}{3} \frac{\epsilon_{r,0}}{kT} = \frac{4}{3} \frac{g_{*S} a_S T_0^4}{2 kT} = 2.890 \times 10^9 \text{ m}^{-3}. \quad (8.13)$$

Recall [from the text immediately after Equation (6.67)] that the (T_ν/T) dependence of g_{*S} is a power of three rather than a power of four, so the factor $(\frac{4}{11})^{4/3}$ becomes just $\frac{4}{11}$ and g_{*S} becomes 3.91.

The present number density of CMB photons is given directly by Equation (6.12):

$$N_\gamma = \zeta(3) \frac{2}{\pi^2} \left(\frac{kT}{c\hbar} \right)^3 = 4.11 \times 10^8 \text{ photons m}^{-3}. \quad (8.14)$$

Neutrino Number Density. Now that we know T_0 and N_γ we can obtain the neutrino temperature $T_\nu = 1.949 \text{ K}$ from Equation (6.65) and the neutrino number density per neutrino species from Equation (6.66),

$$N_\nu = \frac{3}{11} N_\gamma = 1.12 \times 10^8 \text{ neutrinos m}^{-3}. \quad (8.15)$$

For three species of relic neutrinos with average mass $\langle m_\nu \rangle$, Equation (6.68) can be used to cast the density parameter in the form

$$\Omega_\nu = \frac{3 \langle m_\nu \rangle}{94.0 h^2 \text{ eV}}. \quad (8.16)$$

8.2 Temperature Anisotropies

Dipole Anisotropy. The temperature measurement of Penzias and Wilson’s antenna was not very precise by today’s standards. Their conclusion about the isotropy of the CMB was based on an accuracy of only 1.0 K. When the measurements improved over the years it was found that the CMB exhibited a *dipole anisotropy*. The temperature varies minutely over the sky in such a way that it is maximally blueshifted in one direction (call it α) and maximally redshifted in the opposite direction ($\alpha + 180^\circ$). In a direction $\alpha + \theta$ it is

$$T(\theta) = T(\alpha) (1 + v \cos \theta), \quad (8.17)$$

where v is the amplitude of the dipole anisotropy. Although this shift is small, only $vT(\alpha) \approx 3.35$ mK, it was measured with an accuracy better than 1% by the Differential Microwave Radiometer (DMR) instrument on board the COBE satellite [4].

At the end of Chapter 5 we concluded that the hot Big Bang cosmology predicted that the CMB should be essentially isotropic, since it originated in the LSS, which has now receded to a redshift of $z \approx 1080$ in all directions. Note that the most distant astronomical objects known have redshifts of about $z = 7$. Their distance in time to the LSS is actually much closer than their distance to us.

In the standard model the expansion is spherically symmetric, so it is quite clear that the dipole anisotropy cannot be of cosmological origin. Rather, it is well explained by our motion ‘against’ the radiation in the direction of maximal blueshift with relative velocity v .

Thus there is a frame in which the CMB is isotropic—not a rest frame, since radiation cannot be at rest. This frame is then comoving with the expansion of the Universe. We referred to it in Section 2.2, where we noted that, to a fundamental observer at rest in the comoving frame, the Universe must appear isotropic if it is homogeneous. Although general relativity was constructed to be explicitly frame independent, the comoving frame in which the CMB is isotropic is observationally convenient. The fundamental observer is at position B in Figure 8.2.

The interpretation today is that not only does the Earth move around the Sun, and the Solar System participates in the rotation of the Galaxy, but also the Galaxy moves relative to our Local Galaxy Group, which in turn is falling towards a center behind the Hydra–Centaurus supercluster in the constellation Virgo. From the observation that our motion relative to the CMB is 369 ± 0.9 km s^{−1}, these velocity vectors add up to a peculiar motion of the Galaxy of about 550 km s^{−1}, and a peculiar motion of the Local Group of about 630 km s^{−1} [5]. Thus the dipole anisotropy seen by the Earth-based observer A in Figure 8.1 tells us that we and the Local Group are part of a larger, gravitationally bound system.

Multipole Analysis. Temperature fluctuations around a mean temperature T_0 in a direction α on the sky can be analyzed in terms of the *autocorrelation function* $C(\theta)$, which measures the product of temperatures in two directions \mathbf{m}, \mathbf{n} separated by an angle θ and averaged over all directions α ,

$$C(\theta) = \left\langle \frac{\delta T(\mathbf{m})}{T_0} \frac{\delta T(\mathbf{n})}{T_0} \right\rangle, \quad \mathbf{m} \cdot \mathbf{n} = \cos \theta. \quad (8.18)$$

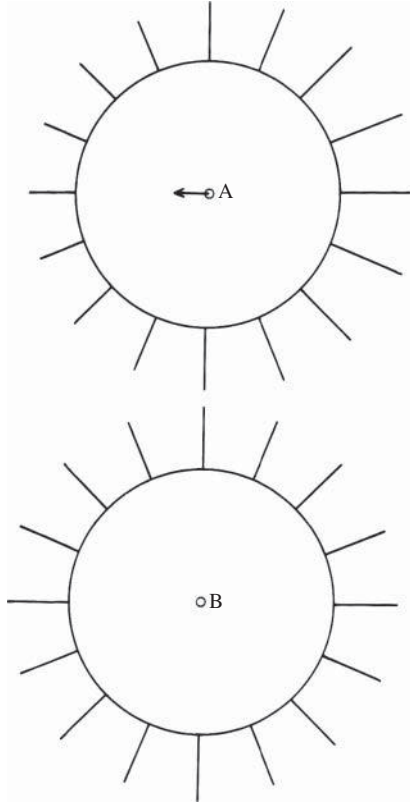


Figure 8.2 The observer A in the solar rest frame sees the CMB to have dipole anisotropy—the length of the radial lines illustrate the CMB intensity—because he is moving in the direction of the arrow. The fundamental observer at position B has removed the anisotropy.

For small angles (θ) the temperature autocorrelation function can be expressed as a sum of *Legendre polynomials* $P_\ell(\theta)$ of order ℓ , the *wavenumber*, with coefficients or *powers* a_ℓ^2 ,

$$C(\theta) = \frac{1}{4\pi} \sum_{\ell=2}^{\infty} a_\ell^2 (2\ell + 1) P_\ell(\cos \theta). \quad (8.19)$$

All analyses start with the quadrupole mode $\ell = 2$ because the $\ell = 0$ monopole mode is just the mean temperature over the observed part of the sky, and the $\ell = 1$ mode is the dipole anisotropy. As a rule of thumb, higher multipoles correspond to fluctuations on angular scales

$$\theta \simeq \frac{60^\circ}{\ell}.$$

In the analysis, the powers a_ℓ^2 are adjusted to give a best fit of $C(\theta)$ to the observed temperature. The resulting distribution of a_ℓ^2 values versus ℓ is called the *power spectrum* of the fluctuations. The higher the angular resolution, the more terms of

high ℓ must be included. Anisotropies on the largest angular scales corresponding to quadrupoles are manifestations of truly primordial gravitational structures.

For the analysis of temperature perturbations over large angles, the Legendre polynomial expansion in Equation (8.19) will not do; one has to use tensor spherical harmonics. The temperature $T(\mathbf{n})$ in the direction \mathbf{n} can be expanded in spherical harmonics

$$T(\mathbf{n}) = T_0 + \sum_{\ell=1}^{\infty} \sum_{m=-\ell}^{\ell} a_{\ell m}^{\text{T}} Y_{\ell m}(\mathbf{n}), \quad (8.20)$$

where $a_{\ell m}^{\text{T}}$ are the powers or *temperature multipole components*. These can be determined from the observed temperature $T(\mathbf{n})$ using the orthonormality properties of the spherical harmonics,

$$a_{\ell m}^{\text{T}} = \frac{1}{T_0} \int d\mathbf{n} T(\mathbf{n}) Y_{\ell m}^*(\mathbf{n}). \quad (8.21)$$

Expressing the autocorrelation function C as a power spectrum in terms of the multipole components, the average of all statistical realizations of the distribution becomes

$$\langle a_{\ell m}^{\text{T}*} a_{\ell' m'}^{\text{T}} \rangle = C_{\ell}^{\text{T}} \delta_{\ell \ell'} \delta_{mm'} = C_{\ell}^{\text{T}}. \quad (8.22)$$

The last step follows from statistical isotropy which requires statistical independence of each lm mode, as manifested by the presence of the two Kronecker deltas and the absence of an m -dependence.

Sources of Anisotropies. Let us now follow the fate of the scalar density perturbations generated during inflation, which subsequently froze and disappeared outside the (relatively slowly expanding) horizon. For wavelengths exceeding the horizon, the distinction between curvature (adiabatic) and isocurvature (isothermal) perturbations is important. Curvature perturbations are true energy density fluctuations or fluctuations in the local value of the spatial curvature. These can be produced, for example, by the quantum fluctuations that are blown up by inflation. By the equivalence principle all components of the energy density (matter, radiation) are affected. Isocurvature fluctuations, on the other hand, are not true fluctuations in the energy density but are characterized by fluctuations in the form of the local equation of state, for example, spatial fluctuations in the number of some particle species. These can be produced, for example, by cosmic strings and other cosmic defects that perturb the local equation of state. As long as an isocurvature mode is superhorizon, physical processes cannot re-distribute the energy density.

When the Universe arrived at the radiation- and matter-dominated epochs, the Hubble expansion of the horizon reveals these perturbations. Once inside the horizon, the crests and troughs can again communicate, setting up a pattern of standing acoustic waves in the baryon-photon fluid. The tight coupling between radiation and matter density causes the adiabatic perturbations to oscillate in phase. After decoupling, the perturbations in the radiation field no longer oscillate, and the remaining standing acoustic waves are visible today as perturbations to the mean CMB temperature at degree angular scales.

Curvature and isocurvature fluctuations behave differently when they are super-horizon: isocurvature perturbations cannot grow, while curvature perturbations can. Once an isocurvature mode passes within the horizon, however, local pressure can move energy density and can convert an isocurvature fluctuation into a true energy-density perturbation. For subhorizon modes the distinction becomes unimportant and the Newtonian analysis applies to both. However, isocurvature fluctuations do not lead to the observed acoustic oscillations seen in Figure 8.3 (they do not peak in the right place), whereas the adiabatic picture is well confirmed.

At the LSS, crests in the matter density waves imply higher gravitational potential. As we learned in Section 2.5, photons ‘climbing out’ of overdense regions will be redshifted by an amount given by Equation (3.1), but this is partially offset by the higher radiation temperature in them. This source of anisotropy is called the *Sachs–Wolfe effect*. Inversely, photons emitted from regions of low density ‘roll down’ from the gravitational potential, and are blueshifted. In the long passage to us they may traverse further regions of gravitational fluctuations, but then their frequency shift upon entering the potential is compensated for by an opposite frequency shift when leaving it (unless the Hubble expansion causes the potential to change during the traverse). They also suffer a time dilation, so one effectively sees them at a different time than unshifted photons. Thus the CMB photons preserve a ‘memory’ of the density fluctuations at emission, manifested today as temperature variations at large angular scales. An anisotropy of the CMB of the order of $\delta T/T \approx 10^{-5}$ is, by the Sachs–Wolfe effect, related to a mass perturbation of the order of $\delta \approx 10^{-4}$ when averaged within one Hubble radius.

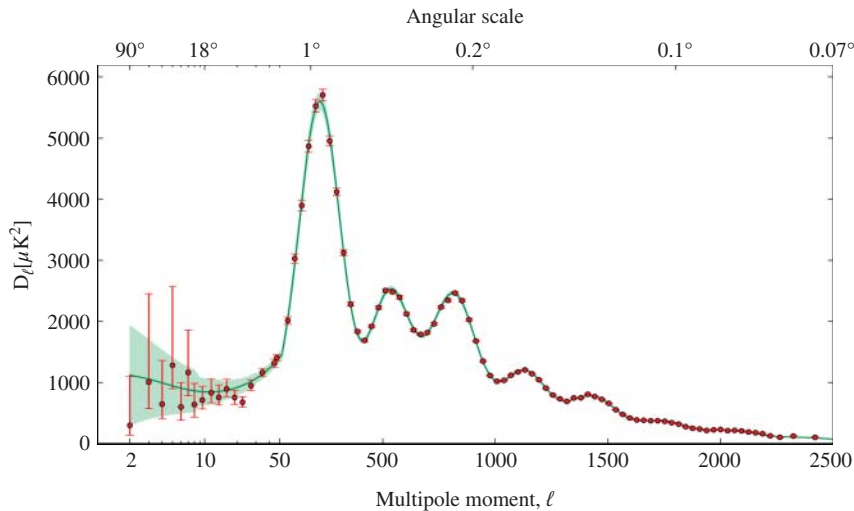


Figure 8.3 The best-fit power spectra of CMB temperature (T) fluctuations as a function of angular scale (top x axis) and multipole moment (bottom x axis) [6]. Reproduced from the freely accessible Planck Legacy Archive with permission of Jan Tauber, European Space Agency. (See plate section for color version.)

The gravitational redshift and the time dilation both contribute to $\delta T/T_0$ by amounts which are linearly dependent on the density fluctuations $\delta\rho/\rho$, so the net effect is given by

$$\frac{\delta T}{T} \simeq \frac{1}{3} \left(\frac{L_{\text{dec}}}{ct_{\text{dec}}} \right)^2 \frac{\delta\rho}{\rho}, \quad (8.23)$$

where L_{dec} is the size of the structure at decoupling time t_{dec} [corresponding to z_{dec} in Equation (5.76)]. [Note that Equation (8.23) is strictly true only for a critical universe with zero cosmological constant.]

The space-time today may also be influenced by primordial fluctuations in the metric tensor. These would have propagated as gravitational waves, causing anisotropies in the microwave background and affecting the large-scale structures in the Universe. High-resolution measurements of the large-angle microwave anisotropy are expected to be able to resolve the tensor component from the scalar component and thereby shed light on our inflationary past.

Further sources of anisotropies may be due to variations in the values of cosmological parameters, such as the cosmological constant, the form of the quintessence potential, and local variations in the time of occurrence of the LSS.

Discovery. For many years microwave experiments tried to detect temperature variations on angular scales ranging from a few arc minutes to tens of degrees. Ever increasing sensitivities had brought down the limits on $\delta T/T$ to near 10^{-5} without finding any evidence for anisotropy until 1992. At that time, the first COBE observations of large-scale CMB anisotropies bore witness of the spatial distribution of inhomogeneities in the Universe on comoving scales ranging from a few hundred Mpc up to the present horizon size, without the complications of cosmologically recent evolution. This is inaccessible to any other astronomical observations.

On board the COBE satellite there were several instruments, of which one, the DMR, received at three frequencies and had two antennas with 7° opening angles directed 60° apart. This instrument compared the signals from the two antennas, and it was sensitive to anisotropies on large angular scales, corresponding to multipoles $\ell < 30$. Later radio telescopes were sensitive to higher multipoles, so one now has a detailed knowledge of the multipole spectrum up to $\ell = 2800$.

The most precise recent results are shown in Figure 8.3. At low ℓ , the temperature–power spectrum is smooth, caused by the Sachs–Wolfe effect. Near $\ell = 200$ it rises towards the first and dominant peak of a series of *Sakharov oscillations*, also confusingly called the *Doppler peak*. They are basically caused by density perturbations which oscillate as acoustic standing waves inside the LSS horizon. The exact form of the power spectrum is very dependent on assumptions about the matter content of the Universe; thus careful measurement of its shape yields precise information about many dynamical parameters. For details of the results included in the figure, see reference [6].

The definitive DMR results [4] cover four years of measurements of eight complete mappings of the full sky followed by the above spherical harmonic analysis. The CMB

anisotropies found correspond to temperature variations of

$$\delta T = 29 \pm 1 \mu\text{K}, \quad \text{or} \quad \delta T/T = 1.06 \times 10^{-5}. \quad (8.24)$$

Half of the above temperature variations, or $\delta T = 15.3 \mu\text{K}$, could be ascribed to quadrupole anisotropy at the 90° angular scale. Although some quadrupole anisotropy is kinetic, related to the dipole anisotropy and the motion of Earth, this term could be subtracted. The remainder is then quadrupole anisotropy of purely cosmological origin.

Since the precision of the COBE measurements surpassed all previous experiments one can well understand that such small temperature variations had not been seen before. The importance of this discovery was succinctly emphasized by the COBE team who wrote that ‘a new branch of astronomy has commenced’. The story of the COBE discoveries have been fascinatingly narrated by George Smoot [7].

8.3 Polarization

The perturbations to the baryon density and the radiation temperature in the tightly coupled baryon–photon fluid are scalar, thus corresponding to a monopole moment ($\ell = 0$). As we saw previously, the radiation field also exhibits dipole perturbations ($\ell = 1$) which are coupled to the baryon bulk velocity, but there are no vector or tensor perturbations. Tensor perturbations would be due to gravitational waves, which have only recently been claimed to have been observed by the BICEP2 Collaboration [8].

The quadrupole moment possessed by free-streaming photons couples more strongly to the bulk velocity (the peculiar velocities) of the baryon–photon fluid than to the density. Therefore, the photon density fluctuations generate temperature fluctuations, while the velocity gradient generates polarization fluctuations.

Thomson Scattering. In Section 6.3 we briefly introduced elastic scattering of photons from free electrons, Equation (6.30), called *Thomson scattering*, or Compton scattering, as it is called for higher frequencies. On that occasion we ignored the fate of the primordial photons, noting only that they were thermalized by this process before decoupling. We also noted that unpolarized photons are polarized by the anisotropic Thomson scattering process, but as long as the photons continued to meet free electrons their polarization was washed out, and no net polarization was produced. At a photon’s last scattering, however, the induced polarization remains and the subsequently free-streaming photon possesses a quadrupole moment ($\ell = 2$).

Consider a plane wave of monochromatic light with frequency ν moving along the momentum vector in the z direction (cf. Figure 8.4). The components of the wave’s electric field vector \mathbf{E} in the (x, y) -plane oscillate with time t in such a way that they can be written

$$E_x(t) = a_x(t) \cos[\nu t - \theta_x(t)], \quad E_y(t) = a_y(t) \cos[\nu t - \theta_y(t)], \quad (8.25)$$

where $a_x(t)$ and $a_y(t)$ are the amplitudes, and $\theta_x(t)$ and $\theta_y(t)$ are the phase angles.

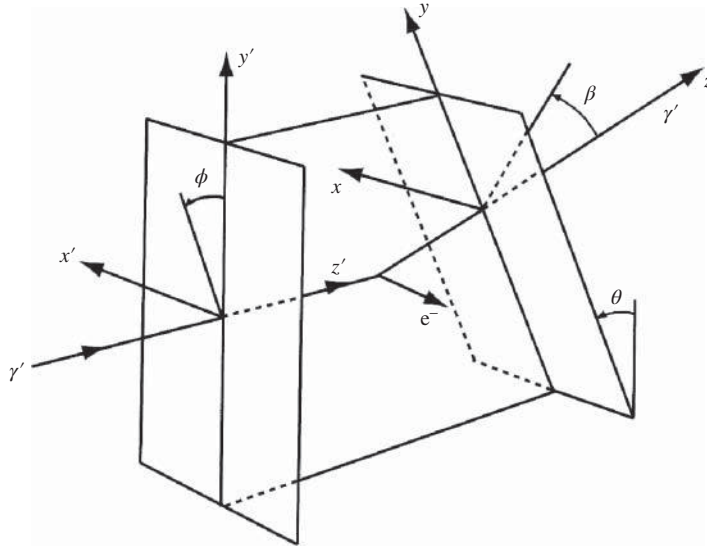


Figure 8.4 The geometry used in the text for describing the polarization of an incoming unpolarized plane wave photon, γ' in the (x', y') -plane, which is Thomson scattering against an electron, and subsequently propagating as a polarized plane wave photon, γ , in the z -direction.

A well-known property of light is its two states of *polarization*. Unpolarized light passing through a pair of polarizing sunglasses becomes vertically polarized. Unpolarized light reflected from a wet street becomes horizontally polarized. The advantage of polarizing sunglasses is that they block horizontally polarized light completely, letting all the vertically polarized light through. Their effect on a beam of unpolarized sunlight is to let, on average, every second photon through vertically polarized, and to block every other photon as if it were horizontally polarized: it is absorbed in the glass. Thus the intensity of light is also reduced to one-half.

Polarized and unpolarized light (or other electromagnetic radiation) can be described by the *Stokes parameters*, which are the time averages (over times much longer than $1/\nu$)

$$\left. \begin{aligned} I &\equiv \langle a_x^2 \rangle + \langle a_y^2 \rangle, & Q &\equiv \langle a_x^2 \rangle - \langle a_y^2 \rangle, \\ U &\equiv \langle 2a_x a_y \cos(\theta_x - \theta_y) \rangle, & V &\equiv \langle 2a_x a_y \sin(\theta_x - \theta_y) \rangle. \end{aligned} \right\} \quad (8.26)$$

The parameter I gives the intensity of light, which is always positive definite. The electromagnetic field is unpolarized if the two components in Equation (8.25) are uncorrelated, which translates into the condition $Q = U = V = 0$. If two components in Equation (8.25) are correlated, they either describe light that is *linearly polarized* along one direction in the (x, y) -plane, or *circularly polarized* in the plane. In the linear case $U = 0$ or $V = 0$, or both. Under a rotation of angle ϕ in the (x, y) -plane, the quantity $Q^2 + U^2$ is an invariant (Problem 2) and the *orientation* of the polarization

$$\alpha \equiv \frac{1}{2} \arctan(U/Q) \quad (8.27)$$

transforms to $\alpha - \phi$. Thus the orientation does not define a direction, it only refers the polarization to the (x, y) -plane.

The photon is peculiar in lacking a longitudinal polarization state, and the polarization is therefore not a vector in the (x, y) -plane; in fact it is a second-rank tensor. This is connected to the fact that the photon is massless. Recall that the theory of special relativity requires the photons to move with the speed of light in any frame. Therefore they must be massless, otherwise one would be able to accelerate them to higher speeds, or decelerate them to rest.

In a way, it appears as if there existed two kinds of photons. Physics has taken this into account by introducing an internal property, *spin*. Thus, one can talk about the two polarization states or about the two spin states of the photon.

Let us now turn to the Stokes parameters in Equation (8.26). The parameter I , which describes the intensity of radiation, is, like V , a physical observable independent of the coordinate system. In contrast, the parameters Q and U depend on the orientation of the coordinate system. In the geometry of Figure 8.4, the coordinates x', y' define a plane wave of incoming radiation propagating in the z' direction (primes are used for unscattered quantities). The incoming photon γ' then Thomson scatters against an electron and the outgoing photon γ continues as a plane wave in a new direction, z .

It follows from the definition of the Stokes parameters Q and U that a rotation of the x' - and y' -axes in the incoming plane by the angle ϕ transforms them into

$$Q(\phi) = Q \cos(2\phi) + U \sin(2\phi), \quad U(\phi) = -Q \sin(2\phi) + U \cos(2\phi). \quad (8.28)$$

We left it as an exercise (Problem 2) to demonstrate that $Q^2 + U^2$ is invariant under the rotation in Equation (8.28). It follows from this invariance that the polarization P is a second rank tensor of the form

$$P = \frac{1}{2} \begin{pmatrix} Q & U - iV \\ U + iV & -Q \end{pmatrix}. \quad (8.29)$$

Thus the polarization is not a vector quantity with a direction unlike the electric field vector \mathbf{E} .

Let us now see how Thomson scattering of the incoming, unpolarized radiation generates linear polarization in the (x, y) -plane of the scattered radiation (we follow closely the pedagogical review of A. Kosowsky [9]). The differential scattering cross-section, defined as the radiated intensity I divided by the incoming intensity I' per unit solid angle Ω and cross-sectional area σ_B , is given by

$$\frac{d\sigma_T}{d\Omega} = \frac{I}{I'} = \frac{3\sigma_T}{8\pi\sigma_B} |i' \cdot i|^2 \equiv K |i' \cdot i|^2. \quad (8.30)$$

Here σ_T is the total Thomson cross-section, the vectors i' , i are the unit vectors in the incoming and scattered radiation planes, respectively (cf. Figure 8.4), and we have lumped all the constants into one constant, K . The Stokes parameters of the outgoing radiation then depend solely on the nonvanishing incoming parameter I' ,

$$I = KI'(1 + \cos^2\theta), \quad Q = KI'\sin^2\theta, \quad U = 0, \quad (8.31)$$

where θ is the scattering angle. By symmetry, Thomson scattering can generate no circular polarization, so $V = 0$ always.

The net polarization produced in the direction \mathbf{z} from an incoming field of intensity $I'(\theta, \phi)$ is determined by integrating Equation (8.31) over all incoming directions. Note that the coordinates for each incoming direction must be rotated some angle ϕ about the z -axis as in Equation (8.31), so that the outgoing Stokes parameters all refer to a common coordinate system. The result is then [9]

$$I(\mathbf{z}) = \frac{1}{2}K \int d\Omega (1 + \cos^2\theta) I'(\theta, \phi), \quad (8.32)$$

$$Q(\mathbf{z}) - iU(\mathbf{z}) = \frac{1}{2}K \int d\Omega \sin^2\theta e^{2i\phi} I'(\theta, \phi). \quad (8.33)$$

Expanding the incident radiation intensity in spherical coordinates,

$$I'(\theta, \phi) = \sum_{\ell m} a_{\ell m} Y_{\ell m}(\theta, \phi), \quad (8.34)$$

leads to the following expressions for the outgoing Stokes parameters:

$$I(\mathbf{z}) = \frac{1}{2}K \left(\frac{8}{3}\sqrt{\pi}a_{00} + \frac{4}{3}\sqrt{\frac{1}{5}}\pi a_{20} \right), \quad (8.35)$$

$$Q(\mathbf{z}) - iU(\mathbf{z}) = 2K\sqrt{\frac{2\pi}{15}}a_{22}. \quad (8.36)$$

Thus, if there is a nonzero quadrupole moment a_{22} in the incoming, unpolarized radiation field, it will generate linear polarization in the scattering plane. To determine the outgoing polarization in some other scattering direction, \mathbf{n} , making an angle β with \mathbf{z} , one expands the incoming field in a coordinate system rotated through β . This derivation requires too much technical detail to be carried out here, so we only state the result [9]:

$$Q(\mathbf{n}) - iU(\mathbf{z}\mathbf{n}) = K\sqrt{\frac{1}{5}}\pi a_{20}\sin^2\beta. \quad (8.37)$$

Multipole Analysis. The tensor harmonic expansion in Equation (8.20) for the radiation temperature T and the temperature multipole components $a_{(\ell m)}^T$ in Equation (8.21) can now be completed with the corresponding expressions for the polarization tensor P . From the expression in Equation (8.29) its components are

$$\begin{aligned} P_{ab}(\mathbf{n}) &= \frac{1}{2} \begin{pmatrix} Q(\mathbf{n}) & -U(\mathbf{n})\sin\theta \\ -U(\mathbf{n})\sin\theta & -Q(\mathbf{n})\sin^2\theta \end{pmatrix} \\ &= T_0 \sum_{\ell=2}^{\infty} \sum_{m=-\ell}^{\ell} \left[a_{(\ell m)}^E Y_{(\ell m)ab}^E(\mathbf{n}) + a_{(\ell m)}^B Y_{(\ell m)ab}^B(\mathbf{n}) \right]. \end{aligned} \quad (8.38)$$

The existence of the two modes (superscripted) E and B is due to the fact that the symmetric traceless tensor in Equation (8.38) describing linear polarization is specified by two independent Stokes parameters, Q and U . This situation bears analogy with the electromagnetic vector field, which can be decomposed into the gradient

of a scalar field (E for electric) and the curl of a vector field (B for magnetic). The source of the E-modes is Thomson scattering. The sources of the B-modes are gravitational waves entailing tensor perturbations, and E-modes which have been deformed by gravitational lensing of large-scale structures in the Universe.

In analogy with Equation (8.21), the polarization multipole components are

$$a_{\ell m}^E = \frac{1}{T_0} \int d\mathbf{n} P_{ab}(\mathbf{n}) Y_{\ell m}^{E ab*}(\mathbf{n}), \quad (8.39)$$

$$a_{\ell m}^B = \frac{1}{T_0} \int d\mathbf{n} P_{ab}(\mathbf{n}) Y_{\ell m}^{B ab*}(\mathbf{n}). \quad (8.40)$$

The three sets of multipole moments $a_{\ell m}^T$, $a_{\ell m}^E$ and $a_{\ell m}^B$ fully describe the temperature and polarization map of the sky; thus, they are physical observables. There are then six power spectra in terms of the multipole components: the temperature spectrum in Equation (8.22) and five spectra involving linear polarization. The full set of physical observables is then

$$\left. \begin{aligned} C_{\ell}^T &= \langle a_{\ell m}^{T*} a_{\ell' m'}^T \rangle, & C_{\ell}^E &= \langle a_{\ell m}^{E*} a_{\ell' m'}^E \rangle, & C_{\ell}^B &= \langle a_{\ell m}^{B*} a_{\ell' m'}^B \rangle, \\ C_{\ell}^{TE} &= \langle a_{\ell m}^{T*} a_{\ell' m'}^E \rangle, & C_{\ell}^{TB} &= \langle a_{\ell m}^{T*} a_{\ell' m'}^B \rangle, & C_{\ell}^{EB} &= \langle a_{\ell m}^{E*} a_{\ell' m'}^B \rangle. \end{aligned} \right\} \quad (8.41)$$

For further details on polarization, see Kosowsky [9].

Thus polarization delivers six times more information than temperature alone, but it is much more difficult to observe since the intensity of polarization is so much weaker than that of temperature.

The first observations of the polarization spectra C_{ℓ}^E and C_{ℓ}^{TE} were made by the South Pole-based Degree Angular Scale Interferometer DASI [10] and the WMAP satellite [7]. We refer to the more recent Planck mission results on E-mode polarization and the South Pole-based BICEP2 telescope results on B-mode polarization [8] in the next section.

8.4 Model Testing and Parameter Estimation

In this section we shall quote best estimators of many cosmological parameters with ‘reliable’ statistical and systematic errors. One may pause to reflect whether that indeed is possible. The problem is that different nonstochastic systematic errors coming from differences in instruments, selections, Monte Carlo modeling, restrictions of the parameter space by different priors and so on, can often not be evaluated and, even less, combined. Moreover, we have only one Universe with its own cosmic variance.

Particle physicists also face the problem of how to combine statistical and systematic errors, and they have some receipts for it.

However, the community of astronomers and astrophysicists is not going to accept any average, because it is not in their culture.

The philosophy of their culture is that the best value comes from one best measurement.

The set of preferred parameter values always shift by some fraction of σ depending on the formulations of the theory, and independently of what data sets are included.

The authors of the currently best measurement then end up presenting ten different versions rather than one recommendation because they do not want to offend anyone.

What one may consider an admissible modification of a model is a cultural question to some extent, and probably this is true even in particle physics. However, in cosmology there appears to be a wider diversity of what people consider to be a reasonable model. For instance, the age of the Universe was quoted in Equation (1.22) to have a finite value with a precision of 50 Myr, but in A. Linde's eternal universe and in the bouncing universe the age may well be infinite.

Thus I caution the reader to take recommended values with a grain of salt.

Expansion Time. The so-called timescale test compares the lookback time in Figure 5.2 at redshifts at which galaxies can be observed with t_0 obtained from other cosmochronometers inside our Galaxy, as discussed in Section 1.4. Thus we have to make do with a consistency test. At moderately high redshifts where the Ω_m term dominates and Ω_λ can be neglected, Equation (5.55) can be written

$$H_0 t(z) \approx \frac{2}{3\sqrt{\Omega_m}}(1+z)^{-3/2}. \quad (8.42)$$

Let us multiply the H_0 and t_0 values in Table A.2 and A.6 to obtain a value for the dimensionless quantity

$$H_0 t_0 = 0.95. \quad (8.43)$$

As we already saw in Equation (5.44) this rules out the spatially flat matter-dominated Einstein–de Sitter universe in which $H_0 t_0 < \frac{2}{3}$.

The Magnitude–Redshift Relation. Equation (2.61) relates the apparent magnitude m of a bright source of absolute magnitude M at redshift z to the luminosity distance d_L . We noted in Section 1.4 that the peak brightness of SNe Ia can serve as remarkably precise standard candles visible from very far away; this determines M . Although the magnitude–redshift relation can be used in various contexts, we are only interested in testing cosmology.

The luminosity distance d_L in Equation 2.60 is a function of z and the model-dependent dynamical parameters, primarily Ω_m , Ω_λ and H_0 . The redshift z or the scale a at the emission of light can be measured in the usual way by observing the shift of spectral lines. The supernova lightcurve shape gives supplementary information: in the rest frame of the supernova the time dependence of light emission follows a standard curve, but a supernova at relativistic distances exhibits a broadened light curve due to time dilation.

Let us define a ‘Hubble-constant free’ luminosity distance

$$D_L(z, \Omega_m \equiv H(z)d_L(z, \Omega_m), \quad (8.44)$$

where $H(z)$ is a simplification of $H(t)$ in Equation 5.55. The apparent magnitude m_B corresponding to d_L (B for the effective B-band, a standard blue filter) becomes

$$m_B = M_B - 5 \log H_0 + 25 + 5 \log D_L(z, \Omega_m, \Omega_\lambda). \quad (8.45)$$

This can be fitted to supernova data to determine best values in the space of the matter density parameter, Ω_m , and the cosmological constant density parameter, Ω_λ ,

introduced in Equation 5.20. The results are in agreement with other determinations, but no longer accurate enough to be of interest.

Reionization. Since polarization originated in the LSS when the horizon was about 1.12° of our present horizon, the polarization fluctuations should only be visible in multipoles

$$\ell > \frac{60^\circ}{1.12^\circ} \approx 54.$$

But the CMB experiments also observe a strong signal on large angular scales, $\ell < 10$. This can only be due to reionization later than LSS, when the radiation on its way to us traversed ionized hydrogen clouds, heated in the process of gravitational contraction. The effect of CMB reionization is called the *Sunyaev–Zel’dovich Effect* (SZE) (*Yakov B. Zel’dovich*, 1914–1987). As a consequence of the SZE, the CMB spectrum is distorted, shifting towards higher energy.

From the size of this spectral shift one estimates a value for the Thomson scattering *optical depth* to the effective reionization clouds, τ [6], which is essentially independent of cosmological modeling

$$\tau = 0.09 \pm 0.01, \tag{8.46}$$

but strongly degenerate with n_s . This corresponds to reionization by an early generation of stars. One quotes [6] the redshift at which the Universe is half reionized, $z_{\text{re}} = 11.15$. It could still be that this picture is simplistic, the reionization may have been a complicated process in a clumpy medium, involving several steps at different redshifts.

CMB Parameters. The parameters required to model CMB are some subset of the parameters H_0 , the observational matter densities $\Omega_m h^2$ and $\Omega_b h^2$ (b for baryons), the cosmological constant density Ω_λ , the optical depth τ , the amplitude A of the power spectrum, the scalar power index n_s in Equation (7.59), the ratio of the tensor to scalar power index, $r = n_t/n_s$, the energy variation of the scalar index dn_s/dk , and the linear theory amplitude of fluctuations σ_8 within $8 \text{ Mpc } h^{-1}$ spheres at $z = 0$.

The primordial fluctuations are assumed to be Gaussian random phase, since no evidence to the contrary has been found.

The first acoustic T peak in Figure 8.3 determines the scale ℓ of the time when matter compressed for the first time after t_{dec} . The positions and amplitudes of the peaks and troughs in the temperature (T) power spectrum Figure 8.3 and the temperature–polarization (TE) cross-power spectrum in Figure 8.5 contain a wealth of information on cosmological parameters. The position in ℓ -space is related to the parameters n_s , $\Omega_m h^2$ and $\Omega_b h^2$. The amplitude of the first peak is positively correlated to $\Omega_m h^2$ and the amplitude of the second peak is negatively correlated to $\Omega_b h^2$ but, to evaluate the physical matter densities Ω_m , Ω_b and Ω_c , one needs to know a value for h , which one can take from Section 1.4. Increasing n_s increases the ratio of the second peak to the first peak. At fixed n_s , however, this ratio determines Ω_b/Ω_m . The amplitudes also determine the effective number of neutrino generations N_{eff} and put limits on the total mass of the neutrino generations, $\Omega_\nu h^2$.

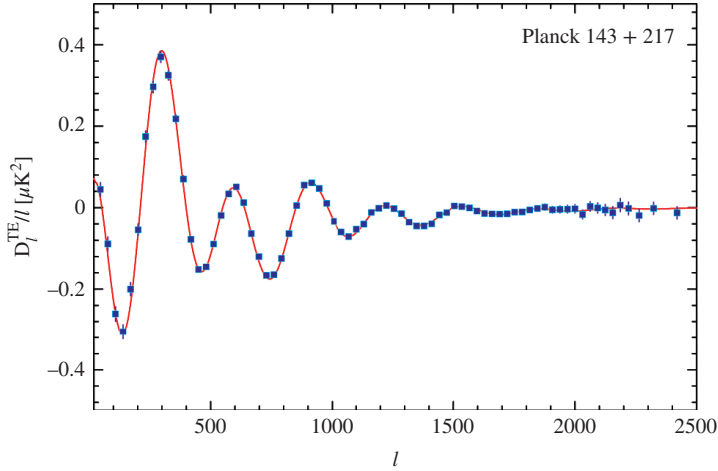


Figure 8.5 The temperature–E-polarization cross-power spectrum as a function of angular scale (top x axis) and multipole moment (bottom x axis) [6]. Reproduced from the freely accessible Planck Legacy Archive with permission of George Efstathiou, Kavli Institute for Cosmology, University of Cambridge. (See plate section for color version.)

Note the qualitative feature in Figure 8.3 that the TE component is zero at the temperature–power spectrum maxima, as is expected (the polarization is maximal at velocity maxima and minimal at temperature maxima, where the velocities are minimal), and that it exhibits a significant large-angle anti-correlation dip at $\ell \approx 150$, a distinctive signature of superhorizon adiabatic fluctuations. The positions and amplitudes of the peaks and troughs in the temperature (T) power spectrum Figure 8.3 and the temperature–polarization (TE) cross-power spectrum in Figure 8.5 contain a wealth of information on cosmological parameters.

In $(\Omega_m, \Omega_\lambda)$ -space, the CMB data determine $\Omega_0 = \Omega_m + \Omega_\lambda$ most precisely, whereas supernova data (discussed in Section 4.4) determine $\Omega_\lambda - \Omega_m$ most precisely. Combining both sets of data with data on large-scale structures from 2dFGRS [11] (discussed in Chapter 9) which depend on $n_s, \Omega_m h, \Omega_b h$ and which put limits on $\Omega_v h$, one breaks the CMB power spectrum parameter degeneracies and improves the precision.

It is too complicated to describe the simultaneous fit to all data and the evaluation of the parameter values here, so we just quote the results published as *Planck+WP* by the Planck Team [6]. This includes information also from the WMAP 9-year polarization low multipole data [12] Note that all errors quoted refer to single-parameter fits at 68% confidence, marginalized over all the other parameters. If one takes into account all recent parameter determinations, the weighted mean central values and the median statistics central values are in fairly good agreement among themselves and with the Planck values, but there is no absolute criterion to select the true values.

Planck+WP finds that the scalar index is

$$n_s = 0.960 \pm 0.0073. \quad (8.47)$$

The South Pole-based BICEP2 telescope [8] has observed the presence of the primordial C_ℓ^B component in a fit to the theoretical B-mode power spectrum. This is reported as a tensor to scalar ratio

$$r = 0.20_{-0.05}^{+0.07}, \quad (8.48)$$

which disfavors $r = 0$ at a confidence of 7.0σ . This result is so remarkable that it certainly needs confirmation. If confirmed, its most likely interpretation is, that the B-mode polarization signals gravitational waves from the time of inflation, and that inflation indeed has occurred.

Density Parameters. The density parameters from all the large experiments combined are [6]

$$\Omega_m h^2 = 0.1385 \pm 0.0025, \quad \Omega_b h^2 = 0.02205 \pm 0.00028. \quad (8.49)$$

The Hubble constant, H_0 , and the matter density parameter, Ω_0 , are only tightly constrained in the combination $\Omega_m h^2$. In the table of Planck + WP results with 68% limits, one finds that the Universe is consistent with being spatially flat, $\Omega_0 = \Omega_m + \Omega_\lambda = 1$. From the combination of all large experiments [8–15]

$$H_0 = 100h = 69.6 \pm 0.7, \quad \Omega_m = 0.286 \pm 0.008. \quad (8.50)$$

From this one can derive the density parameters for baryons (b) and for non-baryonic dark matter (c)

$$\Omega_b = 0.052 \quad \Omega_c = 0.234, \quad (8.51)$$

where we have not evaluated the errors because of model-dependent correlations.

If the Universe is not exactly flat, the vacuum energy is of the order of $\Omega_k = \Omega_0 - 1 = -0.003 \pm 0.003$.

In Figure 8.6 the confidence regions in Ω_m, Ω_λ space are plotted for the combined PanSTARRS1 supernova data [13], Planck CMB data [6], baryonic oscillation BAO data [14] and H_0 data [6].

It is a remarkable success of the FLRW concordance model that the baryonic density at time 380 kyr as evidenced by the CMB is in excellent agreement with the BBN evidence in Equation (6.99) from about 20 min after the Big Bang. As explained in Section 6.4, the BBN value depends only on the expansion rate and the nuclear reaction cross-sections, and not at all on the details of the FLRW model.

In Equation 6.95 we quoted the value for the ratio Y_4 of ^4He mass to total mass $^1\text{H} + ^4\text{He}$ from BBN data, $Y_4^{BBN} = 0.2565 \pm 0.0060$. A more precise value can be found from CMB data, $Y_4^{CMB} = 0.2477 \pm 0.0001$ [6]. A universe with no helium is now ruled out by CMB at very high confidence—it would produce too much small scale power. This provides test of the BBN epoch.

Current limits on the total neutrino mass $\sum m_\nu$ is of the order of <0.66 eV (confidence level CL=95%), but strongly model dependent. Combining this with Equation (8.16), we obtain the ν mass density parameter

$$\Omega_\nu < 0.012. \quad (8.52)$$

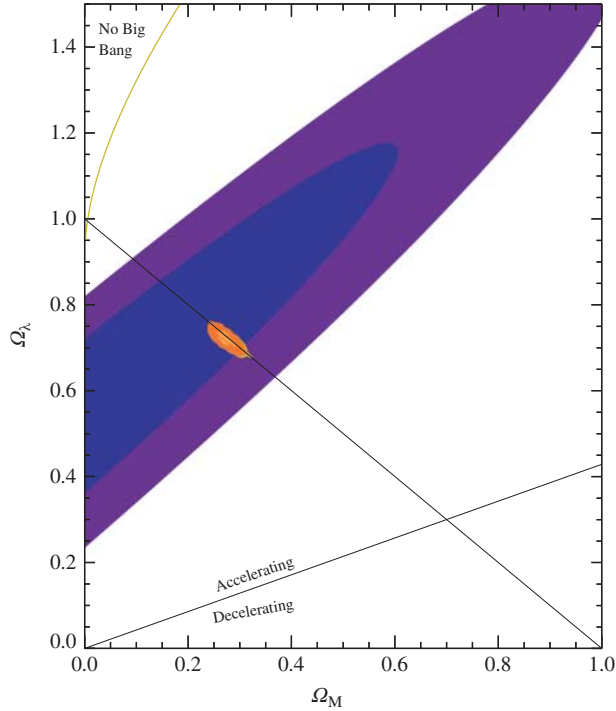


Figure 8.6 The 1σ and 2σ cosmological constraints in Ω_m, Ω_λ space using PanSTARRS1 supernova data [13], Planck CMB data [6], BAO [14] and H_0 data [6] with statistical and systematic errors propagated. Reproduced with permission of Armin Rest for the PanSTARRS1 Collaboration. (See plate section for color version.)

The number of neutrino families is $N_\nu = 3.30 \pm 0.27$ in good agreement with the standard model number 3.

Inserting the values of Ω_b and h in Equation (6.96) we obtain the ratio of baryons to photons

$$\eta = (6.06 \pm 0.08) \times 10^{-10}. \quad (8.53)$$

Inserting the value of h in Equation (8.12), we obtain the present value of the radiation density parameter,

$$\Omega_r = 5.456 \times 10^{-5}. \quad (8.54)$$

Timescales. From this value for Ω_r one can determine the time of equality of radiation and matter density, t_{eq} . From Equations (5.13) and (5.29), the equation determining the evolution of the scale is

$$H(a)^2 = H_0^2[(1 - \Omega_0)a^{-2} + \Omega(a)] = H_0^2[(1 - \Omega_0)a^{-2} + \Omega_m a^{-3} + \Omega_r a^{-4} + \Omega_\lambda].$$

At the present time ($a = 1$) $\Omega_m > \Omega_r$ but, as we move back in time and a gets smaller, the term $\Omega_r a^{-4}$ will come to dominate. The epoch of matter–radiation equality would have

occurred when $\Omega_m a^{-3} = \Omega_r a^{-4}$. Using the value of Ω_m from (8.50) and Ω_r from (8.54) would give $a^{-1} = 1 + z \approx 5773$. This is not correct, however, because the a dependence of Ω_r should actually be given by

$$\Omega_r(a) = \frac{g_*(a) a_S T^4}{2 \rho_c} = \frac{g_*(a) a_S}{2 \rho_c} \left(\frac{2.725}{a} \right)^4, \quad (8.55)$$

using the function g_* discussed in Chapter 5. In the region of $1 + z \gtrsim 1000$ neutrinos will be relativistic and $g_* = 3.36$ instead of 2 [the contribution to the integral in Equation (5.55) from large $z \gtrsim 10^8$, where $g_*(a) > 3.36$ is negligible]. Temperature scales inversely with a , thus at $z_{eq} = 3280$ from WMAP 9-year data [12].

Inserting the values of the energy densities in Equations (8.50) and (8.53) into Equation (5.55), one finds the age of the Universe to a remarkable precision,

$$t_0 = 13.82 \pm 0.05 \text{ Gyr}, \quad (8.56)$$

is in excellent agreement with the independent determinations of lesser precision in Section 1.5. The WMAP team have also derived the redshift and age of the Universe at last scattering and the thickness of the last scattering shell (noting that the WMAP team use the term *decoupling* where we have used *last scattering surface*—see our definitions in Section 6.3):

$$\left. \begin{aligned} t_{\text{LSS}} &= 0.379^{+0.008}_{-0.007} \text{ Myr}, & \Delta t_{\text{LSS}} &= 0.118^{+0.003}_{-0.002} \text{ Myr}, \\ 1 + z_{\text{LSS}} &= 1089 \pm 1, & \Delta z_{\text{LSS}} &= 195 \pm 2. \end{aligned} \right\} \quad (8.57)$$

Deceleration Parameter. Recalling the definitions

$$H_0 = \frac{\dot{a}_0}{a_0}, \quad \Omega_m = \frac{8\pi G \rho_m}{3H_0^2}, \quad \Omega_\lambda = \frac{\lambda}{3H_0^2}, \quad q_0 = -\frac{\ddot{a}_0}{a_0 H_0^2},$$

and ignoring Ω_r , since it is so small, we can find relations between the dynamical parameters Ω_λ , Ω_m , H_0 , and the deceleration parameter q_0 . Substitution of these parameters into Equations (5.17) and (5.18) at present time t_0 , gives

$$H_0^2 + \frac{kc^2}{a_0^2} - \frac{\lambda}{3} = \Omega_m H_0^2, \quad (8.58)$$

$$-2q_0 H_0^2 + H_0^2 + \frac{kc^2}{a_0^2} - \lambda = -3\Omega_m H_0^2 w, \quad (8.59)$$

where w denotes the equation of state $p_m/\rho_m c^2$ of matter.

We can then obtain two useful relations by eliminating either k or λ . In the first case we find

$$\Omega_m(1 + 3w) = 2q_0 + 2\Omega_\lambda, \quad (8.60)$$

and, in the second case,

$$\frac{3}{2}\Omega_m(1 + w) - q_0 - 1 = \frac{kc^2}{S_0^2 H_0^2}. \quad (8.61)$$

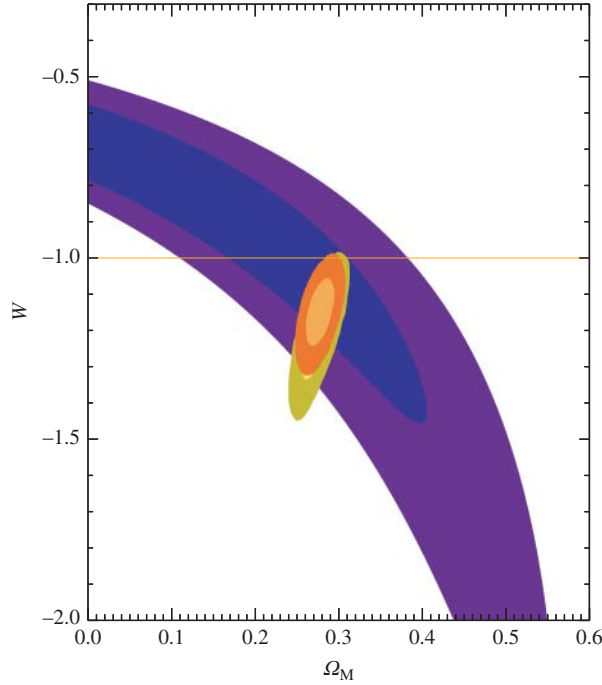


Figure 8.7 The 1σ and 2σ cosmological constraints in Ω_m , W space using PanSTARRS1 supernova data [13], Planck CMB data [6], BAO [14] and H_0 data [6] with statistical and systematic errors propagated. Reproduced with permission of Armin Rest for the PanSTARRS1 Collaboration. (See plate section for color version.)

In the present matter-dominated Universe, the pressure p_m is completely negligible and we can set $w = 0$, as can be seen.

In Figure 8.7 the confidence regions in Ω_m , w space are plotted for the combined PanSTARRS1 supernova data [13], Planck CMB data [6], baryonic oscillation BAO data [14] and H_0 data [6].

From Equation (8.60) and the values for Ω_m and Ω_λ above, we find

$$q_0 = -0.53 \pm 0.02. \quad (8.62)$$

The reason for the very small error is that the errors of Ω_m and Ω_λ are completely anti-correlated. The parameter values given in this section have been used to construct the scales in Figure 6.5. The values are collected in Table A.6 in the Appendix.

Problems

1. Derive an equation for T_{eq} from the condition in Equation (8.3).
2. Show that the quantity $Q^2 + U^2$ is an invariant under the rotation of an angle ϕ in the (x, y) -plane, where Q and U are the Stokes parameters defined in Equation (8.26).

3. Use Wien's constant, Equation (6.120) and the CMB temperature to determine the wavelength of the CMB.
4. Use the present radiation-energy density to calculate the pressure due to radiation in the Universe.
5. Show that an observer moving with velocity β in a direction θ relative to the CMB sees the rest frame blackbody spectrum with temperature T as a blackbody spectrum with temperature

$$T' = \frac{T}{\gamma(1 - \beta \cos \theta)}. \quad (8.63)$$

To first order in β this gives the dipole anisotropy Equation (8.17) [1].

6. The dipole anisotropy is measured to be $1.2 \times 10^{-3}T_0$. Derive the velocity of Earth relative to the comoving coordinate system.
7. Use Equation (5.55) and $\Omega_\lambda = 0.685$ to express t_{eq} in years.

References

- [1] Peebles, P. J. E., Page, Jr., L. A. *et al.* Partridge, R. B. 2009 *Finding the Big Bang cosmology*. Cambridge University Press.
- [2] Mather, J. C., Cheng, E. S., Eplee, R. E. *et al.* 1990 *Astrophys. J. Lett.* **354**, L37.
- [3] Fixsen, D. J. *et al.* 1996 *Astrophys. J.* **473**, 576.
- [4] Bennett, C. L. *et al.* 1996 *Astrophys. J. Lett.* **464**, L1.
- [5] Lynden-Bell, D. *et al.* 1988 *Astrophys. J.* **326**, 19.
- [6] Planck Collaboration: Ade, P. A. R. *et al.* 2014 *Astron. Astrophys* and preprint arXiv:1303.5076 [astro-ph.CO].
- [7] Smoot, G. and Davidson, K. 1993 *Wrinkles in time*. Avon Books, New York.
- [8] BICEP2 Collaboration: Ade, P. A. R. *et al.* 2014 Preprint arXiv:1403.3985 [astro-ph.CO].
- [9] Kosowsky, A. 1999 *New Astronom. Rev.* **43**, 157.
- [10] Kovac, J. *et al.* 2002 *Nature* **420**, 772.
- [11] Colless, M. *et al.* 2001 *Mon. Not. R. Astron. Soc.* **328**, 1039.
- [12] Bennett, C. L. *et al.* 2013 *Astrophys. J. Suppl.* **208**, 20 and earlier WMAP papers.
- [13] Rest, A. *et al.* 2013 Preprint arXiv:1310.3828 [astro-ph.CO].
- [14] Blake, C. *et al.* 2011 *Mon. Not. R. Astron. Soc.* **418**, 1707.
- [15] Bennett, C. L. *et al.* 2014, *Astrophys. J.* **794**, 135.

Dark Matter

The analysis of anisotropies in the cosmic microwave background have shown that the fraction $\Omega_c = 0.263$ of all the matter in the Universe is nonbaryonic. All forms of radiating baryonic mass are already accounted for in Ω_b ; starlight amounts to $\Omega_* = 0.001\text{--}0.002$, gas and star remnants in galaxies amount to $\Omega_{\text{lum}} < 0.01$. The intergalactic gas contains hydrogen clouds and filaments seen by its Ly α absorption, warm gas in groups of galaxies radiates soft X-rays, hot gas in clusters is seen in keV X-rays and in the SZE. The missing fraction is not radiating in any wavelength and is therefore called *dark matter*.

This term is generic for observed gravitational effects on all scales: galaxies, small and large galaxy groups, clusters and superclusters, CMB anisotropies over the full horizon, baryonic oscillations over large scales, and cosmic shear in the large-scale matter distribution. The correct explanation or nature of dark matter is not known, whether it implies unconventional particles or modifications to gravitational theory, but gravitational effects have convincingly proved its existence in some form.

The remarkable fact is that the dark matter fraction is much larger than the known fraction and that it does not interact with ordinary baryonic matter except gravitationally. The purpose of this Chapter is to summarize the phenomenology of all such effects except CMB which was covered in Chapter 8.

In Section 9.1 we study kinematical and dynamical effects: the virial theorem, empirical halo profiles and examples of virially bound groups and clusters. In Section 9.2 we study dark matter in galaxies: spirals, ellipticals, dwarf sphericals, and some arguments on galaxy formation. In Section 9.3 we turn to clusters and in Section 9.4 meet several cases of merging clusters. In Section 9.5 we list possible candidates of dark matter. As we shall see, there are no good candidates, only some hypothetical particles which belong to speculative theories.

In Section 9.6 we turn to observations of galaxy distributions and comparisons with simulations based on the cold dark matter (CDM) paradigm and to predictions and verifications based on the CDM paradigm.

9.1 Virially Bound Systems

The planets move around the Sun along their orbits with orbital velocities balanced by the total gravity of the Solar system. Similarly, stars move in galaxies in orbits with orbital velocities v determined by the gravitational field of the galaxy, or they move with velocity dispersion σ . Galaxies in turn move with velocity dispersion σ under the influence of the gravitational field of their environment, which may be a galaxy group, a cluster or a supercluster. Dark matter forms halos which extend far beyond the luminous matter.

Dynamics. In the simplest dynamical framework one treats massive systems (galaxies, groups and clusters) as statistically steady, spherical, self-gravitating systems of N objects with average mass m and average velocity v or velocity dispersion σ . The total kinetic energy E of such a system is then (we now use σ rather than v)

$$E = (1/2)Nm\sigma^2. \quad (9.1)$$

If the average separation is r , the potential energy of $N(N - 1)/2$ pairings is

$$U = -(1/2)N(N - 1)Gm^2/r. \quad (9.2)$$

The *virial theorem* states that for such a system

$$E = -U/2. \quad (9.3)$$

The total dynamic mass M_{dyn} can then be estimated from σ and r

$$M_{\text{dyn}} = Nm = 2r\sigma^2/G. \quad (9.4)$$

This can also be written

$$\sigma^2 \propto (M_{\text{dyn}}/L)IR, \quad (9.5)$$

where I is a surface luminosity, R is a scale, and M_{dyn}/L is the *mass to light ratio*. Choosing the scale to be the half light radius R_e , this implies a relationship between the observed central velocity dispersion σ_0 , I_e and R_e called the *Fundamental Plane*. of the form

$$R_e \propto (\sigma_0)^a (I_e)^b. \quad (9.6)$$

The virial theorem predicts the values $a = 2$, $b = 1$ for the coefficients. This relationship is found in ellipticals and in some other types of stellar populations with somewhat different coefficients.

Halo Density Profiles. Dark matter halos in galaxies and clusters can of course not be observed by their radiation, but there is an excellent tool to determine their sizes and weights: lensing. This we described in full detail in Section 4.3.

The shapes of DM halos in galaxies and clusters need to be simulated or fitted by empirical formulae. Mostly the shape is taken to be spherically symmetric so that the total gravitating mass profile $M(r)$ depends on three parameters: the mass proportion

in stars, the halo mass and the length scale. A frequently used radial density profile parametrization is

$$\rho_{\text{DM}}(r) = \rho_0 / [(r/r_s)^\alpha (1 + r/r_s)^{3-\alpha}], \quad (9.7)$$

where ρ_0 is a normalization constant and $0 \leq \alpha \leq 3/2$. A standard choice is $\alpha = 1$ for the Navarro–Frenk–White profile (NFW) [1] This presents a singularity at $r \rightarrow 0$ although the total integrated mass is finite. The sharp rise of the density at the halo center forms a “cusp”. Another cusped profile is that of Moore *et al.* [2] with $\alpha = 3/2$.

The Einasto profile [3 and earlier references therein] is defined as

$$\rho_{\text{DM}}(r) = \rho_e \exp\{-d_n[(r/r_e)^{1/n} - 1]\}, \quad (9.8)$$

where the term d_n is a function of n such that ρ_e is the density at r_e , which defines a volume containing half of the total mass. At $r = 0$ the density is then finite and the profile is *cored*.

The Burkert profile [4] has a constant density core

$$\rho_{\text{DM}}(r) = \rho_0 / [(1 + r/r_s)(1 + (r/r_s)^2)], \quad (9.9)$$

which appears to fit dwarf galaxy halos well.

Some clusters are not well fitted by any spherical approximation. The halo may exhibit a strong ellipticity or triaxiality in which case none of the above profiles is good, or the cluster (like Coma) has a binary center.

The dependence of the physical size of clusters on the mass, characterized by the mass concentration index $c \equiv r_{\text{vir}}/r_s$, has been studied in Λ CDM simulations. At intermediate radii c is a crucial quantity in determining the density profile.

The Local Group. The Local Group (LG) is a very small virial system, dominated by two large galaxies, the M31 or Andromeda galaxy, and the Milky Way. The M31 exhibits blueshift, falling in towards us. Evidently our Galaxy and M31 form a bound system together with all or most of the minor galaxies in the Local Group. The Local Group extends to about 3 Mpc and the velocity dispersions of its members is about 200 km s^{-1} . The virial mass is $M_{\text{vir}} = 8.2_{-1.6}^{+2.5} M_{\odot}$.

In this group the two large galaxies dominate the dynamics, so that it is not meaningful to define a statistically average pairwise separation between galaxies, nor an average mass nor an average orbital velocity. The total kinetic energy E is still given by the sum of all the group members, and the potential energy U by the sum of all the galaxy pairs, but here the pair formed by the M31 and the Milky Way dominates, and the pairings of the smaller members with each other are negligible.

Observations of galaxies in the LG can be used to constrain the nature of DM. The number of low mass satellites of MW and M31 as compared with CDM predictions can be explained through the effects of gas-dynamics on baryons, making them invisible, or they can simply not exist, if DM particles have a mass in the keV scale.

Small Galaxy Groups. There are other examples of groups formed by a small number of galaxies which are enveloped in a large cloud of hot gas, visible by its X-ray

emission. One may assume that the electron density distribution associated with the X-ray brightness is in hydrostatic equilibrium, and one can extract the ICM radial density profiles by fits.

The amount of matter in the form of hot gas can be deduced from the intensity of this radiation. Adding the gas mass to the observed luminous matter, the total amount of baryonic matter, M_b , can be estimated. In clusters studied, the gas fraction increases with the distance from the center; the dark matter appears more concentrated than the visible matter.

The temperature of the gas depends on the strength of the gravitational field, from which the total amount of gravitating matter, M_{grav} , in the system can be deduced. In many such small galaxy groups one finds $M_{\text{grav}}/M_b \geq 3$, testifying to a dark halo present. An accurate estimate of M_{grav} requires that also dark energy is taken into account, because it reduces the strength of the gravitational potential. There are sometimes doubts whether all galaxies appearing near these groups are physical members. If not, they will artificially increase the velocity scatter and thus lead to larger virial masses.

On the scale of large clusters of galaxies like the Coma, it is generally observed that DM represents about 85% of the total mass and that the visible matter is mostly in the form of a hot ICM.

The Local Universe. The Local Universe is the best observed part of the universe in which least massive and faintest objects can be detected and studied in detail. These observations resulted in a new research field called Near-Field Cosmology, and have motivated cosmologists to study the Local Group archaeology in their quest for understanding galaxy formation and the play dark matter has on it.

In a volume of a diameter of 96 Mpc beyond the local group Karachentsev [6] has studied 11,000 galaxies appearing single, in pairs, in triplets and in groups. Most of them belong to the Local Supercluster (LSC) and they constitute < 15% of the mass of the Virgo Supercluster. The radial velocities are $v < 3500 \text{ km s}^{-1}$. These galaxies can be treated as a virial system with average density $\Omega_{\text{m,local}} = 0.08 \pm 0.02$, surprisingly small compared to the global density parameter. The fact that much more low mass DM halos are predicted by cosmological simulations than low luminosity galaxies are observed can be explained by gas-dynamical processes which prevent star formation in low mass halos.

Karachentsev quotes [6] three proposed explanations for this mass deficit.

- Dark matter in the systems of galaxies extends far beyond their virial radius, so that the total mass of a group or cluster is three to four times larger than the virial estimate. However, this contradicts other existing data.
- The diameter of the considered region of the Local universe, 90 Mpc, does not correspond to the true scale of the “homogeneity cell”; our Galaxy may be located inside a giant void sized about 100–500 Mpc, where the mean density of matter is three to four times lower than the global value. However, the location of our Galaxy is characterized by an excess, rather than by a deficiency of local density at all scales up to 45 Mpc.

- Most of the dark matter in the Universe, or about two thirds of it, is not associated with groups and clusters of galaxies, but distributed in the space between them in the form of massive dark clumps or as a smooth “ocean”. It is as yet difficult to evaluate this proposal.

Clearly the physics in the Local Universe does not prove the existence of dark matter, rather it brings in new problems.

9.2 Galaxies

Spiral Galaxies. The spiral galaxies are stable gravitationally bound systems in which matter is composed of stars and interstellar gas. Most of the observable matter is in a relatively thin disc, where stars and gas travel around the galactic center on nearly circular orbits. Visible starlight traces velocity out to radial distances typically of the order of 10 kpc, and interstellar gas out to 20–50 kpc. The luminous parts of galaxies, as evidenced by radiation of baryonic matter in the visible, infrared and X-ray spectra, account only for $\Omega_{\text{lum}} < 0.01$.

By observing the Doppler shift of the integrated starlight and the radiation at $\lambda = 0.21$ m from the interstellar hydrogen gas, one finds that spiral galaxies rotate. If the circular velocity at radius r is v in a galaxy of mass M , the condition for stability is that the centrifugal acceleration should equal the gravitational pull:

$$\frac{v^2}{r} = \frac{GM}{r^2}. \quad (9.10)$$

In other words, the radial dependence of the velocity of matter rotating in a disc is expected to follow Kepler’s law $v = \sqrt{GM/r}$.

The surprising result for spiral galaxy rotation curves is, that the velocity does not follow Kepler’s inverse-root law, but stays rather constant after attaining a maximum at about 5 kpc. The most obvious solution to this is that the galaxies are embedded in extensive, diffuse halos of dark matter. In fact, to explain the observations that $v(r) \approx$ constant the radial mass distribution $M(r)$ must be proportional to r and the radial density profile is

$$\rho(r) \propto r^{-2}. \quad (9.11)$$

Assuming that the disc-surface brightness is proportional to the surface density of luminous matter, one derives a circular speed which is typically more than a factor of three lower than the speed of the outermost measured points [5]. This implies that the calculated gravitational field is too small by a factor of 10 to account for the observed rotation.

There are only a few possible solutions to this problem. One is that the theory of gravitation is wrong. It is possible to modify ad hoc Kepler’s inverse square law or Newton’s assumption that G is a constant, but the corresponding modifications cannot be carried out in the relativistic theory, and a general correlation between mass and light remains. The modifications would have to be strong at large scales, and this would greatly enhance cosmic shear, which is inconsistent with measurements.

Another possibility is that spiral galaxies have magnetic fields extending out to regions of tens of kiloparsecs where the interstellar gas density is low and the gas

dynamics may easily be modified by such fields [6]. But this argument works only on the gas halo, and does not affect the velocity distribution of stars. Also, the existence of magnetic fields of sufficient strength remains to be demonstrated; in our Galaxy it is only a few microgauss, which is insufficient.

The accepted solution is then that there exist vast amounts of nonluminous DM beyond that accounted for by luminous, baryonic matter. One natural place to look for DM is in the neighborhood of the Solar System. In 1922, *Jacobus C. Kapteyn* deduced that the total density in the local neighborhood is about twice as large as the luminous density in visible stars. Although the result is somewhat dependent on how large one chooses this ‘neighborhood’ to be, modern dynamical estimates are similar.

Our Galaxy is complicated because of what appears to be a noticeable density dip at 9 kpc and a smaller dip at 3 kpc. To fit the measured rotation curve one needs at least three contributing components: a central bulge, the star disk + gas, and a DM halo. No DM component appears to be needed until radii beyond 15 kpc.

The rotation curve of most galaxies can be fitted by the superposition of contributions from the stellar and gaseous disks, sometimes a bulge, and the dark halo, modeled by a quasi-isothermal sphere. The inner part is difficult to model because the density of stars is high, rendering observations of individual star velocities difficult. Thus the fits are not unique, the relative contributions of disk and dark matter halo is model-dependent, and it is sometimes not even sure whether galactic disks do contain dark matter. Typically, dark matter constitutes about half of the total mass.

In Figure 9.1 we show the rotation curves fitted for 11 well-measured galaxies [7] of increasing halo mass. One notes, that the central dark halo component is indeed much smaller than the luminous disk component. At large radii, however, the need for a DM halo is obvious. On galactic scales, the contribution of DM generally dominates the total mass. Note the contribution of the baryonic component, negligible for light masses but increasingly important in the larger structures.

It appears that cusped profiles are in clear conflict with data on spiral galaxies. Central densities are rather flat, scaling approximately as $\rho_0 \propto r_{\text{luminous}}^{-2/3}$. The best-fit disk + NFW halo mass model fits the rotation curves poorly, it implies an implausibly low stellar mass to light ratio and an unphysically high halo mass. Clearly the actual profiles are of very uncertain origin.

One notes in Figure 9.1 that the shape of the rotation curve depends on the halo virial mass so that the distribution of gravitating matter is luminosity dependent. The old idea that the rotation curve stays constant after attaining a maximum is thus a simplification of the real situation. The rotation velocity can be expressed by a *Universal Rotation Curve* [7]: all spiral galaxies appear to lie on a curve in the four-dimensional space of luminosity, core radius, halo central density and fraction of DM.

What is required to explain the Universal Rotation Curve and the cored profiles is some kind of interaction between baryons and dark matter, which has not met with any success. A more successful idea may be dark matter self-interaction.

Elliptical Galaxies. Elliptical galaxies are quite compact objects which mostly lack neutral gas and which do not rotate so their mass cannot be derived from rotation curves. The total dynamical mass is then the virial mass as derived from the velocity

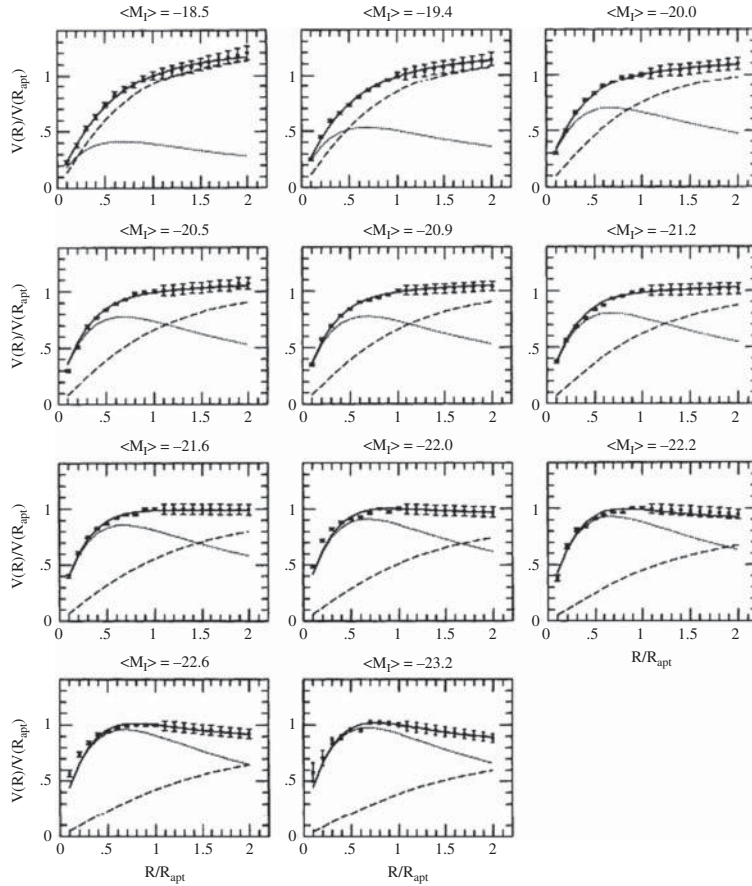


Figure 9.1 The rotation curves fitted for 11 well-measured galaxies of increasing halo mass [7].

dispersions of stars and the anisotropies of their orbits. However, to disentangle the total mass profile into its dark and its stellar components is not straightforward, because the dynamical mass decomposition of dispersions is not unique. The luminous matter in the form of visible stars is a crucial quantity, indispensable to infer the dark component. When available one also makes use of strong and weak lensing data, and of the X-ray properties of the emitting hot gas. The gravity is then balanced by pressure gradients as given by Jeans' Equation (see Chapter 11).

Inside the half light radius R_e the contribution of the dark matter halo to the central velocity dispersion is often very small, it is dominated by the stars, so that the dark matter profile is intrinsically unresolvable. On the average the dark matter component contributes less than 5% to the total velocity dispersions. The outer mass profile is compatible with the NFW Equation (9.7) and Burkert Equation (9.9) formulæ. Important information on the mass distribution can be obtained from the Fundamental Plane, Equation (9.6), which yields the coefficients $a = 1.8$, $b = 0.8$. Note that this

is in some tension with the Virial Theorem, perhaps due to variations in the central dispersions, σ_0 , of the stellar populations.

Most of the baryons in optically selected, isolated, elliptical galaxies are in a morphologically relaxed hot gas halo detectable out to ≈ 200 kpc, that is well described by hydrostatic models. The isolation condition reduces the influence of a possible group-scale or cluster-scale halo. The baryons and the dark matter conspire to produce a total mass density profile that can be well-approximated by a power law, $\rho_{tot} \propto r^{-\alpha}$ over a wide range.

The fitting method involves solving the equation of hydrostatic equilibrium to compute temperature and density profile models, given parametrized mass and entropy profiles. The models are then projected onto the sky and fitted to the projected temperature and density profiles. Fits fit ignoring DM are poor, but the inclusion of DM may improve the fits highly significantly; in one study DM was required at 8.2σ . In several studies, for most of the radii, the dark matter contribution is very small although statistically significant.

Dwarf Spheroidal Galaxies. The mass to light ratio of an astronomical object is defined as $Y \equiv M/L$. Dwarf spheroidal galaxies (dSph) are the smallest stellar systems containing dark matter and exhibit very high $Y = M/L$ ratios, $Y = 10\text{--}100$. In Andromeda IX $Y = 93 + 120/-50$, in Draco $Y = 330 \pm 125$. The dwarf spheroidals have radii of ≈ 100 pc and central velocity dispersions ≈ 10 km s $^{-1}$ which is larger than expected for self-gravitating, equilibrium stellar populations. The generally accepted picture has been, that dwarf galaxies have slowly rising rotation curves and are dominated by dark matter at all radii.

However, Swaters *et al.* [8] have reported observations of H I rotation curves for a sample of 73 dwarf galaxies, among which eight galaxies have sufficiently extended rotation curves to permit reliable determination of the core radius and the central density. They found that dark matter only becomes important at radii larger than three or four disk scale lengths. Their conclusion is, that the stellar disk can explain the mass distribution over the optical parts of the galaxy. Some of the required stellar mass to light ratios are high, up to 15 in the R-band.

Comparing the properties of dwarf galaxies in both the core and outskirts of the Perseus Cluster, Penny and Conselice [9] found a clear correlation between mass to light ratio and the luminosity of the dwarfs, such that the faintest dwarfs require the largest fractions of dark matter to remain bound. This is to be expected, as the fainter a galaxy is, the less luminous mass it will contain, therefore the higher its dark matter content must be to prevent its disruption. Dwarfs are more easily influenced by their environment than more massive galaxies.

The distance to the Perseus Cluster prevents an easy determination of Y , so Penny and Conselice [9] instead determined the dark matter content of the dwarfs by calculating the minimum mass needed in order to prevent tidal disruption by the cluster potential, using their sizes, the projected distance from the cluster center to each dwarf and the mass of the cluster interior. Three of 15 dwarfs turned out to have mass to light ratios smaller than 3, indicating that they do not require dark matter.

Ultra-compact dwarf galaxies (UCDs) are stellar systems with masses of around 10^7 – $10^8 M_{\text{sun}}$ and half-mass radii of 10–100 pc. A remarkable properties of UCDs is that their dynamical mass to light ratios are on average about twice as large as those of globular clusters of comparable metallicity, and also tend to be larger than what one would expect based on simple stellar evolution models. UCDs appear to contain very little or no dark matter.

Collisional N-body simulations of the coevolution of a system composed of stars and dark matter find that DM gets removed from the central regions of such systems due to dynamical friction and mass segregation of stars. The friction timescale is significantly shorter than a Hubble time for typical globular clusters, while most UCDs have friction times much longer than a Hubble time. Therefore, a significant dark matter fraction may remain within the half-mass radius of present-day UCDs, making dark matter a viable explanation for their elevated mass to light ratios.

A different type of systems are the ultra-faint dwarf galaxies (UFDs). When interpreted as steady state objects in virial equilibrium they would be the most DM dominated objects known in the Universe. Their half-light radii range from 70 to 320 pc.

A special case is the UFD disk galaxy *Segue 1* [10] which has a baryon mass of only about 1000 solar masses. One interpretation is that this is a thin non-rotating stellar disk not accompanied by a gas disk, embedded in an axisymmetric DM halo and with a ratio $f \equiv M_{\text{halo}}/M_{\text{b}} \approx 200$. But if the disk rotates, f could be as high as 2000. If *Segue 1* also has a magnetized gas disk, the dark matter halo has to confine the effective pressure in the stellar disk and the magnetic Lorentz force in the gas disk as well as possible rotation. Then f could be very large [10]. Another interpretation is that *Segue 1* is an extended globular cluster rather than an UFD.

Primordial Density Fluctuations. Galaxies form by gas cooling and condensing into DM haloes, where they turn into stars. The star-formation rate is $10M_{\odot} \text{ yr}^{-1}$ in galaxies at $2.8 < z < 3.5$ for which the Ly α break at 91.2 nm shifts significantly (at $z = 3$ it has shifted to 364.8 nm). Galaxy mergers and feedback processes also play major roles.

Galaxy formation requires the study of how galaxies populate DM halos. In simulations one attempts to track galaxy and DM halo evolution across cosmic time in a physically consistent way, providing positions, velocities, star formation histories, abundance matching arguments and other physical properties for the galaxy populations of interest.

The implied spatial clustering of stellar mass turns out to be in remarkably good agreement with a direct and precise measurement. By comparing the galaxy mass autocorrelation function with the mass autocorrelation function averaged over the Local Supercluster (LSC) volume, one concludes that a large amount of matter in the LSC is dark.

Over a wide range of scales there appears to be a universal relation between density and size of observed dark matter halos, from dwarf galaxies to galaxy clusters. Such a universal property is difficult to explain without dark matter.

The very general question arises whether the galaxies could at all have formed from primordial density fluctuations in a purely baryonic medium. As we have noted, the fluctuations in radiation and matter maintain adiabaticity. The amplitude of the primordial baryon density fluctuations would have needed to be very large in order to form the observed number of galaxies. But then the amplitude of the CMB fluctuations would also have been very large, leading to intolerably large CMB anisotropies today. Thus galaxy formation in purely baryonic matter is ruled out by this argument alone. The galaxies could only have formed in the presence of gravitating dark matter which started to fluctuate early, unhindered by radiation pressure.

9.3 Clusters

The Coma Cluster. Historically, the first observation of dark matter in an object at a cosmological distance was made by Fritz Zwicky in 1933 [5]. While measuring radial velocity dispersions of member galaxies in the Coma cluster (that contains some 1000 galaxies), and the cluster radius from the volume they occupy, Zwicky was the first to use the virial theorem to infer the existence of unseen matter. He found to his surprise that the dispersions were almost a factor of ten larger than expected from the summed mass of all visually observed galaxies in the Coma. He concluded that in order to hold galaxies together the cluster must contain huge amounts of some nonluminous matter. From the dispersions he concluded that the average mass of galaxies within the cluster was about 160 times greater than expected from their luminosity (a value revised today), and he proposed that most of the missing matter was dark.

Zwicky's suggestion was not taken seriously at first by the astronomical community which Zwicky felt as hostile and prejudicial. Clearly, there was no candidate for the dark matter because gas radiating X-rays and dust radiating in the infrared could not yet be observed, and nonbaryonic matter was unthinkable—even the neutron had not been discovered yet. Only some 40 years later when studies of motions of stars within galaxies also implied the presence of a large halo of unseen matter extending beyond the visible stars, dark matter became a serious possibility.

Since that time, modern observations have revised our understanding of the composition of clusters. Luminous stars represent a very small fraction of a cluster mass; in addition there is a baryonic, hot *intracluster medium* (ICM) visible in the X-ray spectrum. Rich clusters typically have more mass in hot gas than in stars; in the largest virial systems like the Coma the composition is about 85% DM, 14% ICM and only 1% stars.

In modern applications of the virial theorem one also needs to model and parametrize the radial distributions of the ICM and the dark matter densities. In the outskirts of galaxy clusters the virial radius roughly separates bound galaxies from galaxies which may either be infalling or unbound. The virial radius r_{vir} is conventionally defined as the radius within which the mean density is 200 times the background density.

Matter accretion is in general quite well described within the approximation of the *Spherical Collapse Model*. According to this model, the velocity of the infall motion

and the matter overdensity are related. Mass profile estimation is thus possible once the infall pattern of galaxies is known.

Dark matter is usually dissected from baryons in lensing analyses by first fitting the lensing features to obtain a map of the total matter distribution and then subtracting the gas mass fraction as inferred from X-ray observations. The total mass map can then be obtained with parametric models in which the contribution from cluster-sized DM halos is considered together with the main galactic DM halos. Mass in stars and in stellar remnants is estimated converting galaxy luminosity to mass assuming suitable stellar mass to light ratios.

One may go one step further by exploiting a parametric model which has three kinds of components: cluster-sized DM halos, galaxy-sized (dark plus stellar) matter halos, and a cluster-sized gas distribution. In systems of merging clusters DM may become spatially segregated from baryonic matter and thus observable. We shall meet several such cases later in this Chapter.

The Local Supercluster (LSC). The autocorrelation function $\xi(r)$ in Equation (9.8) was defined for distances r in real space. In practice, distances to galaxies are measured in redshifts, and then two important distortions enter. To describe the separation of galaxy pairs on the surface of the sky, let us introduce the coordinate σ , transversal to the line of sight, and π radial. In redshift space the correlation function is then described by $\xi(\sigma, \pi)$ or its spherical average $\xi(s)$, where $s = \sqrt{\pi^2 + \sigma^2}$.

The transversal distance σ is always accurate, but the radial redshift distance π is affected by velocities other than the isotropic Hubble flow. For relatively nearby galaxies, $r \leq 2$ Mpc, the random peculiar velocities make an unknown contribution to π so that $\xi(s)$ is radially distorted. The undistorted correlation function $\xi(r)$ is seen isotropic in (σ, π) -space in the top left panel of Figure 9.2. The lower left panel of Figure 9.2 shows the distortion to $\xi(s)$ as an elongation in the π direction.

Over large distances where the peculiar velocities are unimportant relative to the Hubble flow (tens of Mpc), the galaxies in the LSC feel its attraction, as is manifested by their infall toward its center with velocities in the range 150–450 km s⁻¹. From this one can derive the local gravitational field and the mass excess δM concentrated in the LSC. The infall velocities cause another distortion to $\xi(s)$: a flattening as is shown in the top right panel of Figure 9.2. When both distortions are included, the correlation function in (σ, π) -space looks like the bottom right panel of Figure 9.2. The narrow peaks in the π direction have been seen for a long time, and are called *Fingers of God*.

If galaxy formation is a local process, then on large scales galaxies must trace mass (on small scales galaxies are less clustered than mass), so that $\xi_{\text{gal}}(r)$ and $\xi_{\text{mass}}(r)$ are proportional:

$$\xi_{\text{gal}}(r) = b^2 \xi_{\text{mass}}(r).$$

Here b is the linear *bias*: bias is when galaxies are more clustered than mass, and *anti-bias* is the opposite case; $b = 1$ corresponds to the unbiased case. The presence of bias is an inevitable consequence of the nonlinear nature of galaxy formation. The distortions in $\xi(s)$ clearly depend on the mass density Ω_m within the observed volume.

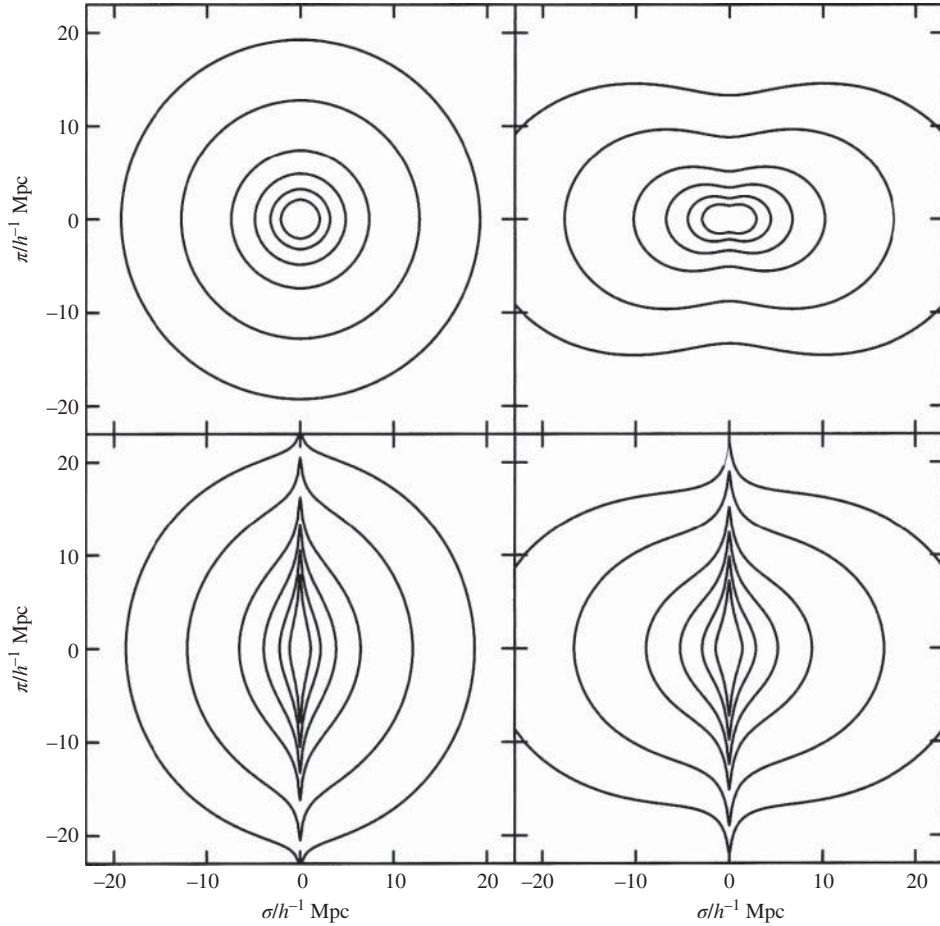


Figure 9.2 Plot of theoretically calculated correlation functions $\xi(\sigma, \pi)$ as described in reference [2]. The lines represent contours of constant $\xi(\sigma, \pi) = 4.0, 2.0, 0.5, 0.2, 0.1$. The models are: top left, undistorted; bottom left, no infall velocities but $\beta = 0.4$; top right, infall velocity dispersion $a = 500 \text{ km s}^{-1}$, $\beta = 0$; bottom right, $a = 500 \text{ km s}^{-1}$, $\beta = 0.4$. All models use Equation (9.15) with $r_c = 5.0 h^{-1} \text{ Mpc}$ and $\gamma = 1.7$. Reproduced from [11] by permission of the 2dFGRS Team.

Introducing the phenomenological flattening parameter

$$\beta = \Omega_m^{0.6} / b, \quad (9.12)$$

one can write a linear approximation to the distortion as

$$\frac{\xi(s)}{\xi(r)} = 1 + \frac{2\beta}{3} + \frac{\beta^2}{5}. \quad (9.13)$$

Estimates of β and b which we shall discuss in Equations (9.18–19) lead to a value of Ω_m of the order of 0.25. Thus a large amount of matter in the LSC is dark.

9.4 Merging Galaxy Clusters

In isolated galaxies and galaxy clusters all matter components contributing to the common gravitational potential are more or less centrally-symmetrically coincident. This makes it difficult to discern DM from the baryonic components, and dependent on parametrization, as we have discussed in the section on Dynamics. In merging galaxy clusters however, the separate distributions of galaxies, intracluster gas and DM may become spatially segregated permitting separate observations. The visually observable galaxies behave as collisionless particles, the baryonic intracluster plasma is fluid-like, it experiences ram pressure and emits X-rays, but noninteracting DM does not feel that pressure, it only makes itself felt by its contribution to the common gravitational potential.

Major cluster mergers are the most energetic events in the Universe since the Big Bang. Shock fronts in the intracluster gas are the key observational tools in the study of these systems. When a subcluster traverses a larger cluster it cannot be treated as a solid body with constant mass moving at constant velocity. During its passage through the gravitational potential of the main cluster it is shrinking over time, stripped of gas envelope and decelerating. Depending on the ratio of the cluster masses, the gas forms a bow shock in front of the main cluster, and this can even be reversed at the time when the potentials coincide.

In several examples of galaxy cluster mergers the presence of DM could be inferred from the separation of the gravitational potential from the position of the radiating plasma.

The Bullet Cluster 1E0657-558. The exceptionally hot and X-ray luminous galaxy cluster 1E0657-558, the *Bullet cluster* at redshift $z = 0.296$, was discovered by Tucker et al. in 1995 [13] in *Chandra* X-ray data. Its structure as a merger of a $2.3 \times 10^{14} M_{\odot}$ subcluster with a main $2.8 \times 10^{14} M_{\odot}$ cluster was demonstrated as the first clear example of a bow shock in a heated intracluster plasma. With the advent of high-resolution lensing, a technique could be developed combining multiple strongly-lensed *Hubble Space Telescope* multi-color images of identified galaxies with weakly lensed and elliptically distorted background sources. The reconstructed gravitational potential does not trace the X-ray plasma distribution which is the dominant baryonic mass component, but rather approximately traces the distribution of bright cluster member galaxies.

In early pictures, often reproduced, the center of the total mass was offset from the center of the baryonic mass peaks, proving that the majority of the matter in the system is unseen. In front of the bullet cluster which has traversed the larger one about 100 Myr ago with a relative velocity of 4500 km s^{-1} , a bow shock is evident in the X-rays. The main cluster peak and the distinct subcluster mass concentration are both clearly offset from the location of the X-ray gas. We do not reproduce that well known figure here because recent analysis of this system [14] do not confirm these results, rather they find that dark matter forms three distinct clumps.

The Baby Bullet. Another merging system with similar characteristics but with lower spatial resolution is the post-merging galaxy cluster pair MACS J0025.4-1222, also called the *Baby Bullet*. It has an apparently simple geometry, consisting of two large subclusters of similar richness, about $2.5 \times 10^{14} M_{\odot}$, both at redshift $z = 0.586$, colliding in approximately the plane of the sky. Multiple images due to strong lensing of four distinct components could be identified.

The two distinct mass peaks are clearly offset by 4σ from the main baryonic component, which is the radiating hot gas observed by *Chandra*. The relative merging velocity is estimated to be 2000 km s^{-1} . The majority of the mass is spatially coincident with the identified galaxies, which implies that the cluster must be dominated by a relatively collisionless form of dark matter.

Many clusters exhibit very complicated geometries observed by *Chandra*. However, in many cases the separation of hot gas from collisionless DM still requires better lensing data.

“El Gordo.” The Atacama Cosmology Telescope has presented properties for an exceptionally massive merging cluster, the ACT-CL J0102-4915 nicknamed *El Gordo* at redshift $z = 0.87$ [15]. It was discovered by selection for its bright Sunyaev-Zel’dovich (SZ) effect (described in Section 10.2), confirmed optically and through its *Chandra* X-ray data. It is the most significant SZ cluster detection to date by nearly a factor of two, with an SZ decrement comparable to the Bullet cluster.

The galaxy distribution is double peaked, whereas the peak in the X-ray emission lies between the density peaks. The X-ray peak forms a relatively cool bullet of low entropy gas like in the Bullet cluster. In the absence of a weak lensing mass reconstruction, the galaxy distribution can only be used as a proxy for the total mass distribution. Thus to conclude that an offset between baryonic and DM has been demonstrated is yet premature.

Other Mergers. Data acquired with the *Advanced Camera for Surveys* on the *Hubble Space Telescope*, *HST*, shows that the cluster Abell 2744 [16] is a complicated merger between three or four separate bodies. The position and mass distribution of part of the cluster have been tightly constrained by the strong lensing of 11 background galaxies producing 31 multiple images. But two clumps lack such images from strong lensing, indicating that they are less massive.

The joint gravitational lensing analysis combines all the strongly lensed multiply-imaged systems and their redshifts with weak lensing shear catalogues to reconstruct the cluster’s lensing potential. The location of shock fronts and velocities, densities and temperatures in the intracluster medium, and all existing X-ray data from *Chandra* have been combined. The lensing mass reconstruction and the luminosity contours of the emission then shows an extremely complex picture of separations between the dark matter and baryonic components. However, it is difficult to ascertain whether this is a single, separate DM structure and to derive decisive separation between DM, X-ray luminous gas and bright cluster member galaxies. One possible interpretation is a near simultaneous double merger 0.12–0.15 Gyr ago.

Another major cluster merger is DLSC J0916.2+2951 [17] at $z = 0.53$, in which the collisional cluster gas has become clearly dissociated from the collisionless galaxies and dark matter. The cluster was identified using optical and weak-lensing observations as part of the Deep Lens Survey. Follow-up observations with Keck, Subaru, Hubble Space Telescopes, and Chandra show that the cluster is a dissociative merger which constrains the DM self-interaction cross-section to $\sigma m_{\text{DM}}^{-1} \leq 7 \text{ cm}^2/\text{g}$. The system is observed at least $0.7 \pm 0.2 \text{ Gyr}$ since first pass-through, thus providing a picture of cluster mergers 2–5 times further progressed than similar systems observed to date.

9.5 Dark Matter Candidates

If only a small percentage of the total mass of the Universe is accounted for by stars and hydrogen clouds, could baryonic matter in other forms make up DM? The answer given by nucleosynthesis is a qualified no: all baryonic DM is already included in Ω_{b} .

Dark Baryonic Matter. Before the value of Ω_{dm} was pinned down as certainly as it is now, several forms of dark baryonic matter was considered. Gas or dust clouds were the first thing that came to mind. We have already accounted for hot gas because it is radiating and therefore visible. Clouds of cold gas would be dark but they would not stay cold forever. Unless there exists vastly more cold gas than hot gas, which seems unreasonable, this DM candidate is insufficient.

It is known that starlight is sometimes obscured by *dust*, which in itself is invisible if it is cold and does not radiate. However, dust grains re-radiate starlight in the infrared, so they do leave a trace of their existence. But the amount of dust and rocks needed as DM would be so vast that it would have affected the composition of the stars. For instance, it would have prevented the formation of low-metallicity (population-II) stars. Thus dust is not an acceptable candidate. And baryons are just too few to explain DM.

Snowballs of frozen hydrogenic matter, typical of comets, have also been considered, but they would sublimate with time and become gas clouds. A similar fate excludes *collapsed stars*: they eject gas which would be detectable if their number density were sufficient for DM.

A more serious candidate for baryonic matter has been *jupiters* or *brown dwarfs*: stars of mass less than $0.08 M_{\odot}$. They also go under the acronym MACHO for *Massive Compact Halo Object*. They lack sufficient pressure to start hydrogen burning, so their only source of luminous energy is the gravitational energy lost during slow contraction. Such stars would clearly be very difficult to see since they do not radiate. However, if a MACHO passes exactly in front of a distant star, the MACHO would act as a gravitational microlens, because light from the star bends around the massive object. The intensity of starlight would then be momentarily amplified (on a timescale of weeks or a few months) by microlensing, as described in Section 4.3. The problem is that even if MACHOs were relatively common, one has to monitor millions of stars for one positive piece of evidence. Only a few microlensing MACHOs have been discovered in the space between Earth and the Large Magellanic Cloud, but not enough to explain DM.

The shocking conclusion is that the predominant form of matter in the Universe is nonbaryonic, and we do not even know what it is composed of! Thus we are ourselves made of some minor pollutant, a discovery which may well be called the fourth breakdown of the anthropocentric view. The first three were already accounted for in Chapter 1.

Black Holes. Primordial black holes could be good candidates because they evade the nucleosynthesis bound, they are not luminous, they (almost) do not radiate, and if they are big enough they have a long lifetime, as we saw in Equation (5.82). They are believed to sit at the center of every galaxy and have masses exceeding $100M_{\odot}$. The mass range $0.3\text{--}30M_{\odot}$ is excluded by the nonobservation of MACHOs in the galactic halo (cf. Section 5.4). Various astrophysical considerations limit their mass to around 10^4M_{\odot} . But black holes cannot be the solution to the galactic rotation curves nor to merging clusters.

Exotica. At higher scales extra dimensions could exist. the usual(3+1) spacetime could be a *brane* embedded in a higher dimensional *bulk*. Standard model fields would be confined on the brane while gravity propagates in the extra dimension. In this kind of models the hierarchy problem is solved in various ways for example the higher dimensions are compactified on different topologies of scale R with the effect of lowering the Planck scale energy closer to the electroweak scale. Compactification of extra dimensions gives rise to a quantization of momentum, $p^2 \simeq 1/R^2$ of the fields propagating in the bulk, and the apparition of a set of Fourier expanded modes [Kaluza Klein (KK) states], for each bulk field. Particles moving in extra dimensions appear as heavy particles with masses m_n/R . The new states have the same quantum numbers (e.g., charge, color, etc). Another way to reach the same goal is to introduce extra dimensions with large curvature.

If extra dimensions are compactified around a circle or torus, the extra-dimensional momentum conservation implies conservation of the KK number n , and the lightest first level KK state is stable. Theories with compact extra dimensions can be written as theories in ordinary four dimensions by performing a KK reduction.

Among further exotica are *solitons*, which are nontopological scalar-field quanta with conserved global charge Q (*Q-balls*) or baryonic charge B (*B-balls*).

Mirror Matter [18] implies reintroducing, all the known fields with the same coupling constants, but with opposite parities. The most natural way to do so is to add to the existing Lagrangian its parity-symmetric counterpart, so that the whole Lagrangian is invariant under the parity transformation, each part transforming into the other. We then end up with a new sector of particles, called mirror sector, which is an exact duplicate of the ordinary sector, but where ordinary particles have left-handed interactions, mirror particles have right-handed interactions.

As a consequence, the three gauge interactions act separately in each sector, the only link between them being gravity. Because mirror baryons, just like their ordinary counterparts, are stable and can be felt only through their gravitational effects, the mirror matter scenario provides an ideal interpretation of dark matter. Its particularity is that it is a self-interacting candidate, but without any new parameter at the

level of particle physics. At the cosmological level there are two new parameters, the ratio x of the temperatures of the ordinary and mirror cosmic background radiations, and the relative amount β of mirror baryons compared to the ordinary ones. In the presence of mirror matter, the matter energy density parameter is

$$\Omega_m = \Omega_b(1 + \beta). \quad (9.14)$$

Scalar fields Suppose that the dark matter particles are described by a spin-0 scalar field like the QCD axion. A variety of unification theories have proposed other scalar field candidates, the bosonic particles being typically ultralight with an ultrahigh phase-space density. This may lead to a *Bose–Einstein Condensate* (BEC), a macroscopic occupancy of the many-body ground state [19]. In principle, for a fixed number of thermalized identical bosons, a BEC will form if a $n\lambda_{\text{deB}}^3 \gg 1$ where n is the number density and λ_{deB} is the de Broglie wavelength. This is equivalent to there being a critical temperature T_c , below which a BEC can form. For small boson masses the corresponding critical temperature of condensation is so high ($\gg \text{TeV}$), that the bosons are fully condensed very early on, that is, almost all of the bosons occupy the lowest available energy state. Hence, the cosmological scalar field dark matter can be described by a single coherent, classical scalar field.

If the thermal decoupling within the bosonic dark matter occurs when the expansion rate exceeds its thermalization rate well after condensation, most of the bosons will stay in the ground state (BEC) and the classical field remains a good description, analogous to the fact that CMB photons after decoupling still follow a black-body distribution. Given the huge critical temperature at hand, one may effectively consider the BEC state as an initial condition. On the other hand, one may also envisage a scenario in which the coherent scalar field is created gravitationally at the end of inflation. A prime motivation for studying scalar field dark matter has been its ability to suppress small-scale clustering and hence potentially resolve some dark matter problems.

Supersymmetric Cold Dark Matter (CDM). Particles which were very slow at time t_{eq} when galaxy formation started are candidates for CDM. If these particles are massive and have weak interactions, so called WIMPs (*Weakly Interacting Massive Particles*), they became nonrelativistic much earlier than the leptons and become decoupled from the hot plasma. For instance, the supersymmetric models (SUSY) contain a very large number of new particles, of which the lightest ones would be stable. At least three such neutral SUSY ‘sparticles’—the *photino*, the *Zino* and the *Higgsino*—or a linear combination of them (the *neutralino*) could serve. However, negative results from laboratory searches now appear to rule out the minimal supersymmetric model as a remedy for CDM. Heavier nonminimally supersymmetric particles are no longer obvious CDM candidates.

Sterile Neutrinos. Very heavy sterile neutrinos could also be CDM candidates, or other cold thermal relics of mass up to some 300 TeV. All this is very speculative. Alternatively, the CDM particles may be very light if they have some superweak

interactions, in which case they froze out early when their interaction rate became smaller than the expansion rate, or they never even attained thermal equilibrium. Candidates in this category are the *axion* and its SUSY partner *axino*. The axion is a light pseudoscalar boson with a 2γ coupling like the π^0 , so it could convert to a real photon by exchanging a virtual photon with a proton. Its mass is expected to be of the order of $1\ \mu\text{eV}$ to $10\ \text{meV}$. It was invented to prevent CP violation in QCD, and it is related to a slightly broken baryon number symmetry in a five-dimensional space-time. Another CDM candidate could be axion clusters with masses of the order of $10^{-8} M_\odot$.

The WIMPs would traverse terrestrial particle detectors with a typical virial velocity of the order of $200\ \text{km s}^{-1}$, and perhaps leave measurable recoil energies in their elastic scattering with protons. The proof that the recoil detected was due to a particle in the galactic halo would be the annual modulation of the signal. Because of the motion of the Earth around the Sun, the signal should have a maximum in June and a minimum in December. Several experiments to detect such signals are currently running with controversial or no results which only permit setting upper limits to the WIMP flux.

All WIMPs have in common that they are hitherto unobserved particles which only are predicted by some theories. A signal worth looking for would be monoenergetic photons from their annihilation

$$X_{\text{dm}} + \bar{X}_{\text{dm}} \rightarrow 2\gamma. \quad (9.15)$$

No such convincing signals have been observed. All reviews to date only list upper limits.

The supersymmetric WIMP scenario is very uneconomical since it duplicates the standard model with its large number of particles when actually only one new particle would be needed for DM. In this sense also the mirror model is uneconomical.

DM Distribution. The ideal fluid approximation which is true for the collisionless DM on large scales breaks down when they decouple from the plasma and start to stream freely out of overdense regions and into underdense regions, thereby erasing all small inhomogeneities (*Landau damping*). This defines the characteristic length and mass scales for freely streaming particles of mass m_{dm} ,

$$\lambda_{\text{fs}} \simeq 40 \left(\frac{30\ \text{eV}}{m_{\text{dm}}} \right) \text{Mpc}, \quad (9.16)$$

$$M_{\text{fs}} \simeq 3 \times 10^{15} \left(\frac{30\ \text{eV}}{m_{\text{dm}}} \right)^2 M_\odot. \quad (9.17)$$

Perturbations in CDM start growing from the time of matter–radiation equality, while baryonic fluctuations are inhibited until recombination because of the tight coupling with photons (or alternatively one can say because of the large baryonic Jeans mass prior to recombination). After recombination, the baryons fall into the CDM potential wells. A few expansion times later, the baryon perturbations catch up with the DM, and both then grow together until $\delta > 1$, when perturbations become Jeans unstable (cf. Chapter 10), collapse and virialize. The amplitude of radiation, however, is unaffected by this growth, so the CMB anisotropies remain at the level

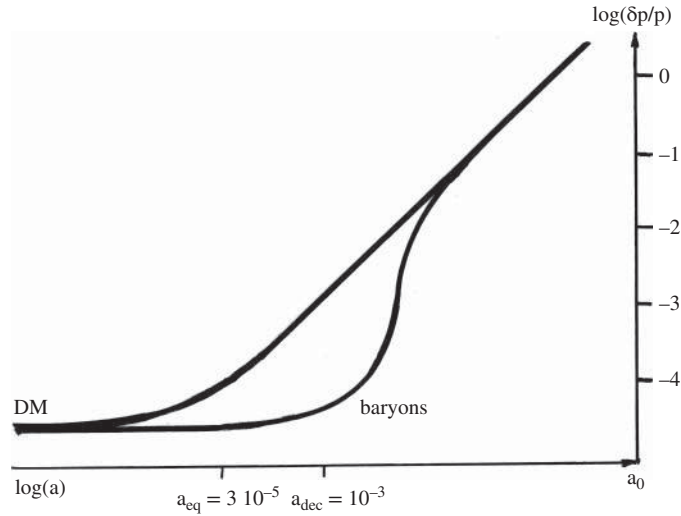


Figure 9.3 Time-dependence of the DM and baryon density fluctuations, $\delta\rho/\rho$. The perturbations in DM start to grow near the epoch of matter–radiation equality, t_{eq} . However, the perturbations in the baryons cannot begin to grow until just after decoupling, when baryons fall into the DM potential wells. Within a few expansion times the baryon perturbations catch up with the DM perturbations.

determined by the baryonic fluctuations just before recombination. This is illustrated in Figure 9.3.

The lightest DM particles are slow enough at time t_{eq} to be bound in perturbations on galactic scales. They should then be found today in galaxy haloes together with possible baryonic MACHOs. If the nonbaryonic DM in our galactic halo were to be constituted by DM at a sufficient density, they should also be found inside the Sun, where they lose energy in elastic collisions with the protons, and ultimately get captured. They could then contribute to the energy transport in the Sun, modifying solar dynamics noticeably. So far this possible effect has not been observed.

On the other hand, if the DM overdensity only constituted early potential wells for the baryons, but did not cluster so strongly, most WIMPs would have leaked out now into the intergalactic space. In that case the DM distribution in clusters would be more uniform than the galaxy (or light) distribution, so that galaxies would not trace mass.

Hot and Warm Dark Matter. Although the neutrinos decoupled from the thermal plasma long before matter domination, they remained relativistic for a long time because of their small mass. For this reason they would possibly constitute *hot dark matter* (HDM), freely streaming at t_{eq} . The accretion of neutrinos to form haloes around the baryon clumps would be a much later process. The CMB is then very little perturbed by the clumps, because most of the energy is in neutrinos and in radiation. However, we already know that the neutrino fraction is much too small to make it a DM candidate, so HDM is no longer a viable alternative.

An intermediate category is constituted by possible sterile neutrinos and by the *gravitino*, which is a SUSY partner of the graviton. These have been called *warm dark matter* (WDM). Both HDM and WDM are now ruled out by computer simulations of the galaxy distribution in the sky. WDM is also ruled out by the CMB detection of early re-ionization at $z > 11$. We shall therefore not discuss these alternatives further.

9.6 The Cold Dark Matter Paradigm

The λ CDM paradigm (also written λ CDM or LCDM) is based on all the knowledge we have assembled so far: the FLRW model with a spatially flat geometry, BBN and thermodynamics with a known matter inventory including dark energy (cf. Chapter 11) of unknown origin but known density; inflation-caused linear, adiabatic, Gaussian mass fluctuations accompanying the CMB anisotropies with a nearly scale-invariant Harrison–Zel’dovich power spectrum; growth by gravitational instability from t_{eq} until recombination, and from hot gas to star formation and hierarchical clustering.

The new element in this scenario is collisionless DM, which caused matter domination to start much earlier than if there had been only baryons. The behavior of DM is governed exclusively by gravity (unless we discover any DM interactions with matter or with itself), whereas the formation of the visible parts of galaxies involves gas dynamics and radiative processes.

While the CMB temperature and polarization anisotropies measure fluctuations at recombination, the galaxy distribution measures fluctuations up to present times. Cosmic shear in weak lensing is sensitive to the distribution of DM directly, but it leaves a much weaker signal than do clusters.

Currently observable bright galaxies are more than 170 billion in number, out to 230 Mpc. The comoving future visibility limit is about 19 Gpc counting how far a photon can travel from the Big Bang to the infinite future. The total number of galaxies that one will eventually be able to see is then 400 billion, and the absolute limit in redshift is $z < 1.69$ [12].

Distributions of galaxies in two-dimensional pictures of the sky show that they form long filaments separating large underdense voids with diameters up to $60 h^{-1}$ Mpc. Figure 9.4 shows such a map out to 2.74 Gyr [20]. The image reveals a wealth of detail, including linear supercluster features, often nearly perpendicular to the line of sight. The largest structure is the Sloan Great Wall of galaxies $1.37 h^{-1}$ billion lightyears long, 80% longer than the previously known Great Wall.

Hierarchical Scenarios. Early CDM models (without an Ω_λ component) produced galaxies naturally, but underproduced galaxy clusters and supergalaxies of mass scale $10^{15} M_\odot$. This was an example of a bottom–top scenario, where small-scale structures were produced first and large-scale structures had to be assembled from them later. Although there was not time enough in this scenario to produce large-scale structures within the known age of the Universe, the scenario could be improved by the introduction of a new degree of freedom, the cosmological constant.

The opposite ‘top–bottom’ scenario was predicted by HDM models where the first structures, supergalaxies, formed by neutrino clouds contracting into pancakes which

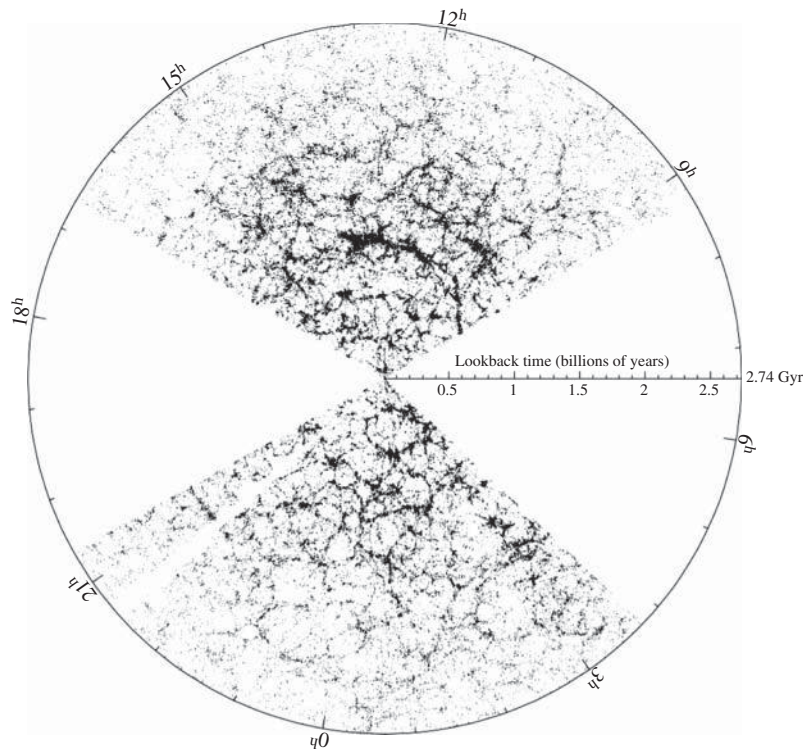


Figure 9.4 Zoom-in of the region of the sky showing galaxies observed by the SDSS Collaboration out to a radius of 2.74 Gyr [20]. The scorpion-like structure in the upper quadrant is the Sloan Great Wall. From Gott III, J. *et al.*, A Map of the Universe, *Astrophys. J.*, **624**, 463, published May 10 2005. © AAS. Reproduced with permission.

subsequently collapsed and disintegrated. Smaller structures and galaxies formed later from the crumbs. But computer simulations of pancake formation and collapse show that the matter at the end of the collapse is so shocked and so heated that the clouds do not condense but remain ionized, unable to form galaxies and attract neutrino haloes. Moreover, large clusters (up to $10^{14} M_{\odot}$) have higher escape velocities, so they should trap five times more neutrinos than large galaxies of size $10^{12} M_{\odot}$. This scenario is not supported by observations, which show that the ratio of dynamic mass to luminous mass is about the same in objects of all sizes, except for dwarf galaxies.

The ‘bottom–top’ scenario is supported by the observations that supergalaxies are typically at distances $z \lesssim 0.5$, whereas the oldest objects known are quasars at redshifts up to $z = 5\text{--}7$. There are also several examples of galaxies which are older than the groups in which they are now found. Moreover, in our neighborhood the galaxies are generally falling in towards the Virgo cluster rather than streaming away from it.

Several pieces of evidence indicate that luminous galaxies could have been assembled from the merging of smaller star-forming systems before $z \approx 1$. The Hubble Space Telescope as well as ground-based telescopes have discovered vast numbers of faint

blue galaxies at $1 \leq z \leq 3.5$, which obviously are very young. There is also evidence that the galaxy merger rate was higher in the past, increasing roughly as a^{-m} or $(1+z)^m$ with $m \approx 2-3$. All this speaks for a bottom-top scenario.

Large Scale Structure Simulation. The formation and evolution of cosmic structures is so complex and nonlinear and the number of galaxies considered so enormous that the theoretical approach must make use of either numerical simulations or semi-analytic modeling. The strategy in both cases is to calculate how density perturbations emerging from the Big Bang turn into visible galaxies. This requires a number of processes in a phenomenological manner:

- (i) the growth of DM haloes by accretion and mergers;
- (ii) the dynamics of cooling gas;
- (iii) the transformation of cold gas into stars;
- (iv) the spectrophotometric evolution of the resulting stellar populations;
- (v) the feedback from star formation and evolution on the properties of prestellar gas;
- (vi) the build-up of large galaxies by mergers.

The primary observational information consists of a count of galaxy pairs in the redshift space coordinates σ , π . From this, the correlation function $\xi(s)$ in redshift space, and subsequently the correlation function $\xi(r)$ in real space, can be evaluated. Here $\xi(s)$ and $\xi(r)$ are related via the parameter β in Equation (9.13). From $\xi(r)$, the power spectrum $P(k)$ can in principle be constructed using its definition later, in Equations (10.8) and (10.9).

The observed count of galaxy pairs is compared with the count estimated from a randomly generated mass distribution following the same selection function both on the sky and in redshift. Different theoretical models generate different simulations, depending on the values of a large number of adjustable parameters: h , $\Omega_m h \equiv (\Omega_{\text{dm}} h^2 + \Omega_b h^2)/h$, Ω_b/Ω_m , Ω_0 , n_s , the normalization σ_8 and the bias b between galaxies and mass.

The CDM paradigm sets well-defined criteria on the real fluctuation spectrum. A good fit then results in parameter values. Since the parameter combinations here are not the same as in the CMB analysis, the degeneracy in the 2dFGRS data between $\Omega_m h$ and Ω_b/Ω_m can be removed by combining the CMB and 2dFGRS analyses. Let us now summarize a few of the results.

If the simulated mass-correlation function $\xi_{\text{dm}}(r)$ and the observed galaxy-number two-point correlation function $\xi_{\text{gal}}(r)$ are identical, this implies that light (from galaxies) traces mass exactly. If not, they are biased to a degree described by b . The result is that there is no bias at large scales, as indeed predicted by theory, but on small scales some anti-bias is observed. This result is a genuine success of the theory because it does not depend on any parameter adjustments. Independently, weak lensing observations also show that visible light in clusters does trace mass (all the visible light is emitted by the stars in galaxies, not by diffuse emission), but it is not clear whether this is true on galaxy scales.

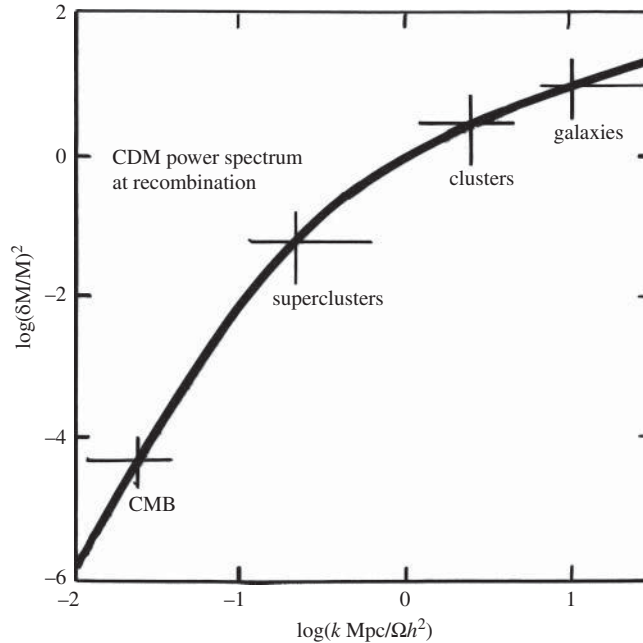


Figure 9.5 The theoretical power function $P(k)$ in the Λ CDM model as a function of the density fluctuation wavenumber k in units of $h^2 \text{Mpc}^{-1}$. This can also be expressed by the angular scale in degrees or by the linear size L of present cosmic structures in units of $\text{Mpc} \Omega_m^{-1} h^{-1}$. The crosses indicate the approximate locations of galaxies, clusters, superclusters and the CMB.

Theoretical models predict that the brightest galaxies at $z = 3$ should be strongly clustered, which is indeed the case. This comparison is also independent of any parameter adjustments. In contrast, DM is much more weakly clustered at $z = 3$ than at $z = 0$, indicating that galaxies were strongly biased at birth.

In Figure 9.5 we show the theoretical linear power spectrum $P(k)$. The real 2dFGRS galaxy power spectrum data lie so accurately on the theoretical curve in Figure 9.5 that we have refrained from plotting them. To achieve this success, all the free parameters have been adjusted.

One important signature of gravitational instability is that material collapsing around overdense regions should exhibit peculiar velocities and infall leading to redshift-space distortions of the correlation as shown in Figure 9.3. We have previously referred to large-scale bulk flows of matter observed within the LSC attributed to the ‘Great Attractor’, an overdensity of mass about $5.4 \times 10^{16} M_\odot$ in the direction of the *Hydra–Centaurus* cluster, but far behind it, at a distance of some $79 h^{-1} \text{Mpc}$. The 2dFGRS has verified that both types of redshift-space distortions occur, the ‘Fingers of God’ due to nearby peculiar velocities and the flattening due to infall at larger distances. These results are quantified by the parameter value

$$\beta = \Omega_m^{0.6} / b = 0.43 \pm 0.07. \quad (9.18)$$

With a large-scale bias of $b = 1$ to a precision of about 10%, the β error dominates, so that one obtains

$$\Omega_m = 0.25 \pm 0.07, \quad (9.19)$$

in good agreement with other determinations.

Later than the above results the Planck collaboration [21] obtained a value for the characteristic amplitude of velocity fluctuations within $8 \text{ Mpc } h^{-1}$ spheres at $z = 0$,

$$\sigma_8 = 0.83 \pm 0.01.$$

To summarize one can state that on scales larger than a few Mpc the distribution of DM in CDM models is essentially understood. Understanding the inner structure of DM haloes and the mechanisms of galaxy formation has proved to be much more difficult.

Problems

1. Suppose that galaxies have flat rotation curves out to R_{max} . The total mass inside R_{max} is given by Equation (9.10), where v may be taken to be 220 km s^{-1} . If the galaxy number density is $n = 0.01 h^3 / \text{Mpc}^3$, show that $\Omega = 1$ when R_{max} is extended out to an average intergalactic distance of $2.5 h^{-1} \text{ Mpc}$ [18].
2. Suppose that neutralinos have a mass of 100 GeV and that they move with a virial velocity of 200 km s^{-1} . How much recoil energy would they impart to a germanium nucleus?

References

- [1] Navarro, J.F. *et al.* 1997 *Astrophys. J.* **490**, 493.
- [2] Moore B. *et al.* 1998 *Astrophys. J. Lett.* **499**, L5.
- [3] Einasto J. *et al.* 1974 *Nature* **250** 309.
- [4] Burkert A. 1995 *Astrophys. J. Lett.* **447**, L25.
- [5] Zwicky, F. 1933 *Helvetica Physica Acta* **6**, Issue 2, 110.
- [6] Karachentsev, I. D. 2012 *Astrophys. Bulletin* **67**, Issue 2, 123.
- [7] Cattaneo, A., Salucci P. and Papastergis, E. 2014 *Astrophys. J.* **783**, 66.
- [8] Swaters, R. A. *et al.* 2011 *Astrophys. J.* **729**, 118.
- [9] Penny, S. J. and Conselice, C. J. 2011 *EAS Publications Series* **48**, 197.
- [10] Xiang-Gruess, M. *et al.* 2009 *Mon. Not. Roy. Astr. Soc.* **400**, L52.
- [11] Hawkins, E. *et al.* 2002 arXiv astro-ph/0212375 and *Mon. Not. R. Astron. Soc.* (2003).
- [12] Ciarcelluti, P. and Wallemacq, Q. 2014 preprint arXiv:1211.5354 [astro-ph.CO].
- [13] Tucker, W. H. *et al.* 1998 *Astrophys. J.* **496**, L5.
- [14] Paraficz, D. *et al.* 2012 preprint arXiv:1209.0384 [astro-ph.CO].
- [15] Menenteau, F. *et al.* 2012 *Astrophys. J.* **748**, 18.
- [16] Merten, J. *et al.* 2011 *Mon. Not. Roy. Astr. Soc.* **417**, 333.
- [17] Dawson, W. E. *et al.* 2012 *Astrophys. J.* **747**, L42.
- [18] Li, B *et al.* 2014 *Phys. Rev. D* and preprint arXiv:1310.6061 [astro-ph.CO].
- [19] Peebles, P. J. E. 1993 *Principles of Physical Cosmology*. Princeton University Press, Princeton, NJ.
- [20] Gott III, J. *et al.* 2005 *Astrophys. J.* **624**, 463.
- [21] Planck Collaboration: Ade, P. A. R. *et al.* 2014 preprint arXiv:1303.5076 [astro-ph.CO] and The Planck Legacy Archive.

Cosmic Structures

After the decoupling of matter and radiation described in Chapter 6, we followed the fate of the free-streaming CMB in Chapter 8. Here we shall turn to the fate of matter and cold nonradiating dust. After recombination, when atoms formed, density perturbations in baryonic matter could start to grow and form structures, but growth in weakly interacting nonbaryonic (dark) matter could have started earlier at the time of radiation and matter equality. The time and size scales are important constraints to galaxy-formation models, as are the observations of curious patterns of filaments, sheets and voids on very large scales.

In Section 10.1 we describe the theory of density fluctuations in a viscous fluid, which approximately describes the hot gravitating plasma. This very much parallels the treatment of the fluctuations in radiation that cause anisotropies in the CMB.

In Section 10.2 we learn how pressure and gravitation conspire so that the hot matter can begin to cluster, ultimately to form the perhaps 10^9 galaxies, clusters and other large-scale structures.

10.1 Density Fluctuations

Until now we have described the dynamics of the Universe by assuming homogeneity and adiabaticity. The homogeneity cannot have grown out of primeval chaos, because a chaotic universe can grow homogeneous only if the initial conditions are incredibly well fine-tuned. Vice versa, a homogeneous universe will grow more chaotic, because the standard model is gravitationally unstable.

But the Universe appears homogeneous only on the largest scales (a debatable issue!), since on smaller scale we observe matter to be distributed in galaxies, groups of galaxies, supergalaxies and strings of supergalaxies with great voids in between. At the time of matter and radiation equality, some lumpiness in the energy density must have been the 'seeds' or *progenitors* of these cosmic structures, and one would expect to see traces of that lumpiness also in the CMB temperature anisotropies originating

in the last scattering. The angular scale subtended by progenitors corresponding to the largest cosmic structures known, of size perhaps $200 h^{-1}$ Mpc, is of the order of 3° , corresponding to CMB multipoles around $\ell = 20$.

Viscous Fluid Approximation. The common approach to the physics of matter in the Universe is by the hydrodynamics of a viscous, nonstatic fluid. With this nonrelativistic (Newtonian) treatment and linear perturbation theory we can extract much of the essential physics while avoiding the necessity of solving the full equations of general relativity. In such a fluid there naturally appear random fluctuations around the mean density $\bar{\rho}(t)$, manifested by compressions in some regions and rarefactions in other regions. An ordinary fluid is dominated by the material pressure but, in the fluid of our Universe, three effects are competing: radiation pressure, gravitational attraction and density dilution due to the Hubble flow. This makes the physics different from ordinary hydrodynamics: regions of overdensity are gravitationally amplified and may, if time permits, grow into large inhomogeneities, depleting adjacent regions of underdensity.

The nonrelativistic dynamics of a compressible fluid under gravity is described by three differential equations, the *Eulerian equations*. Let us denote the density of the fluid by ρ , the pressure p , and the velocity field \mathbf{v} , and use comoving coordinates, thus following the time evolution of a given volume of space. The first equation describes the conservation of mass: what flows out in unit time corresponds to the same decrease of matter in unit space. This is written

$$\frac{d\rho}{dt} = -\rho \nabla \cdot \mathbf{v}. \quad (10.1)$$

Next we have the equation of motion of the volume element under consideration,

$$\frac{d\mathbf{v}}{dt} = -\frac{1}{\rho} \nabla p - \nabla \phi, \quad (10.2)$$

where ϕ is the gravitational potential obeying Poisson's equation, which we met in Equation (3.35),

$$\nabla^2 \phi = 4\pi G \rho. \quad (10.3)$$

Equation (10.2) shows that the velocity field changes when it encounters pressure gradients or gravity gradients.

The description in terms of the Eulerian equations is entirely classical and the gravitational potential is Newtonian. The Hubble flow is entered as a perturbation to the zeroth-order solutions with infinitesimal increments $\delta\mathbf{v}$, $\delta\rho$, δp and $\delta\phi$. Let us denote the local density $\rho(\mathbf{r}, t)$ at comoving spatial coordinate \mathbf{r} and world time t . Then the fractional departure at \mathbf{r} from the spatial mean density $\bar{\rho}(t)$ is the dimensionless *mass density contrast*

$$\delta_m(\mathbf{r}, t) = \frac{\rho_m(\mathbf{r}, t) - \bar{\rho}_m(t)}{\bar{\rho}_m(t)}. \quad (10.4)$$

The solution to Equations (10.1)–(10.3) can then be sought in the form of waves,

$$\delta_m(\mathbf{r}, t) \propto e^{i(\mathbf{k}\cdot\mathbf{r} - \omega t)}, \quad (10.5)$$

where \mathbf{k} is the wave vector in comoving coordinates. An arbitrary pattern of fluctuations can be described mathematically by an infinite sum of independent waves, each with its characteristic wavelength λ or comoving wavenumber k and its amplitude δ_k . The sum can be formally expressed as a Fourier expansion for the density contrast

$$\delta_m(\mathbf{r}, t) \propto \sum \delta_k(t) e^{i\mathbf{k}\cdot\mathbf{r}}. \quad (10.6)$$

A density fluctuation can also be expressed in terms of the mass M moved within one wavelength, or rather within a sphere of radius λ , thus $M \propto \lambda^3$. It follows that the wavenumber or spatial frequency k depends on mass as

$$k = \frac{2\pi}{\lambda} \propto M^{-1/3}. \quad (10.7)$$

Power Spectrum. The density fluctuations can be specified by the amplitudes δ_k of the dimensionless *mass autocorrelation function*

$$\xi(r) = \langle \delta(\mathbf{r}_1) \delta(\mathbf{r} + \mathbf{r}_1) \rangle \propto \sum \langle |\delta_k(t)|^2 \rangle e^{i\mathbf{k}\cdot\mathbf{r}}, \quad (10.8)$$

which measures the correlation between the density contrasts at two points \mathbf{r} and \mathbf{r}_1 . The powers $|\delta_k|^2$ define the power spectrum of the root-mean-squared (RMS) mass fluctuations

$$P(k) = \langle |\delta_k(t)|^2 \rangle. \quad (10.9)$$

Thus the autocorrelation function $\xi(r)$ is the Fourier transform of the power spectrum. We have already met a similar situation in the context of CMB anisotropies, where the waves represented temperature fluctuations on the surface of the surrounding sky. There Equation (8.18) defined the autocorrelation function $C(\theta)$ and the powers a_ℓ^2 were coefficients in the Legendre polynomial expansion Equation (8.19).

Taking the power spectrum to be of the phenomenological form [Equation (7.60)],

$$P(k) \propto k^n,$$

and combining with Equation (10.7), one sees that each mode δ_k is proportional to some power of the characteristic mass enclosed, M^α .

Inflationary models predict that the mass density contrast obeys

$$\delta_m^2 \propto k^3 \langle |\delta_k(t)|^2 \rangle \quad (10.10)$$

and that the primordial fluctuations have approximately a Harrison–Zel’dovich spectrum with $n_s = 1$. Support for these predictions come from the CMB temperature and polarization asymmetry spectra which give the value quoted in Equation (8.47), $n = 0.960 \pm 0.0073$ [1].

Independent, although less accurate information about the spectral index can be derived from constraints set by CMB isotropy, galaxies and black holes using $\delta_k \propto M^\alpha$. The CMB scale (within the size of the present horizon of $M \approx 10^{22} M_\odot$) is isotropic to less than about 10^{-4} . Galaxy formation (scale roughly $10^{12} M_\odot$) requires perturbations of order $10^{-4 \pm 1}$. Taking the ratio of perturbations versus mass implies a constraint on α and implies that n is close to 1.0 at long wavelengths.

Turning to short wavelengths (scale about 10^{12} kg or $10^{-18} M_\odot$), black holes provide another constraint. Primordial perturbations on this scale must have been roughly smaller than 1.0. Larger perturbations would have led to overproduction of black holes, since large-amplitude perturbations inevitably cause overdense regions to collapse before pressure forces can respond. From Equation (5.82) one sees that black holes less massive than 10^{12} kg will have already evaporated within 10 Gyr, but those more massive will remain and those of mass 10^{12} kg will be evaporating and emitting γ rays today. Large amplitude perturbations at and above 10^{12} kg would imply more black holes than is consistent with the mass density of the Universe and the γ ray background. Combining the black hole limit on perturbations (up to around 1 for $M \approx 10^{12}$ kg) with those from the CMB and galaxy formation also implies the spectrum must be close to the Harrison–Zel’dovich form.

The power spectra of theoretical models for density fluctuations can be compared with the real distribution of galaxies and galaxy clusters. Suppose that the galaxy number density in a volume element dV is n_G , then one can define the probability of finding a galaxy in a random element as

$$dP = n_G dV. \quad (10.11)$$

If the galaxies are distributed independently, for instance with a spatially homogeneous Poisson distribution, the joint probability of having one galaxy in each of two random volume elements dV_1, dV_2 is

$$dP_{12} = n_G^2 dV_1 dV_2. \quad (10.12)$$

There is then no correlation between the probabilities in the two elements. However, if the galaxies are clustered on a characteristic length r_c , the probabilities in different elements are no longer independent but correlated. The joint probability of having two galaxies with a relative separation r can then be written

$$dP_{12} = n_G^2 [1 + \xi(r/r_c)] dV_1 dV_2, \quad (10.13)$$

where $\xi(r/r_c)$ is the *two-point correlation function* for the galaxy distribution. This can be compared with the autocorrelation function in Equation (10.8) of the theoretical model. If we choose our own Galaxy at the origin of a spherically symmetric galaxy distribution, we can simplify Equation (10.13) by setting $n_G dV_1 = 1$. The right-hand side then gives the average number of galaxies in a volume element dV_2 at distance r .

Analyses of galaxy clustering show [2] that, for distances

$$10 \text{ kpc} \lesssim hr \lesssim 10 \text{ Mpc}, \quad (10.14)$$

a good empirical form for the two-point correlation function is

$$\xi(r/r_c) = (r/r_c)^{-\gamma}, \quad (10.15)$$

with the parameter values $r_c \approx 5.0 h^{-1} \text{ Mpc}$, $\gamma \approx 1.7$.

Irregularities in the metric can be expressed by the curvature radius r_U defined in Equation (5.54). If r_U is less than the linear dimensions d of the fluctuating region, it will collapse as a black hole. Establishing the relation between the curvature of the metric and the size of the associated mass fluctuation requires the full machinery of general relativity, which is beyond our ambitions.

Linear Approximation. Much of the interesting physics of density fluctuations can be captured by a Newtonian linear perturbation analysis of a viscous fluid. Small perturbations grow slowly over time and follow the background expansion until they become heavy enough to separate from it and to collapse into gravitationally bound systems. As long as these perturbations are small they can be decomposed into Fourier components that develop independently and can be treated separately. For fluctuations in the linear regime, $|\delta_k| < 1$, where

$$\rho_m = \bar{\rho}_m + \Delta\rho_m, \quad p = \bar{p} + \Delta p, \quad v^i = \bar{v}^i + \Delta v^i, \quad \phi = \bar{\phi} + \Delta\phi, \quad (10.16)$$

the size of the fluctuations and the wavelengths grows linearly with the scale a , whereas in the nonlinear regime, $|\delta_k| > 1$, the density fluctuations grow faster, with the power a^3 , at least (but not exponentially). The density contrast can also be expressed in terms of the linear size d of the region of overdensity normalized to the curvature radius,

$$\delta \approx \left(\frac{d}{r_U} \right)^2. \quad (10.17)$$

In the linear regime r_U is large, so the Universe is flat. At the epoch when d is of the order of the Hubble radius, the density contrast is

$$\delta_H \approx \left(\frac{r_H}{r_U} \right)^2, \quad (10.18)$$

free streaming can leave the region and produce the CMB anisotropies. Structures formed when $d \ll r_H$, thus when $\delta \ll 1$. Although δ may be very small, the fluctuations may have grown by a very large factor because they started early on (see Problem 3 in Chapter 7).

When the wavelength is below the horizon, causal physical processes can act and the (Newtonian) viscous fluid approximation is appropriate. When the wavelength is of the order of or larger than the horizon, however, the Newtonian analysis is not sufficient. We must then use general relativity and the choice of gauge is important.

Gauge Problem. The mass density contrast introduced in Equation (9.4) and the power spectrum of mass fluctuations in Equation (10.9) represented perturbations to an idealized world, homogeneous, isotropic, adiabatic, and described by the FLRW model. For subhorizon modes this is adequate. For superhorizon modes one must apply a full general-relativistic analysis. Let us call the space-time of the world just described \mathcal{G} . In the real world, matter is distributed as a smooth background with mass perturbations imposed on it. The space-time of that world is not identical to \mathcal{G} , so let us call it \mathcal{G}' .

To go from \mathcal{G}' , where measurements are made, to \mathcal{G} , where the theories are defined, requires a *gauge transformation*. This is something more than a mere coordinate transformation—it also changes the event in \mathcal{G} that is associated to an event in \mathcal{G}' . A perturbation in a particular observable is, by definition, the difference between its value at some space-time event in \mathcal{G} and its value at the corresponding event in the background (also in \mathcal{G}). An example is the mass autocorrelation function $\xi(r)$ in Equation (10.8).

But this difference need not be the same in \mathcal{G}' . For instance, even if an observable behaves as a scalar under coordinate transformations in \mathcal{G} , its perturbation will not be invariant under gauge transformations if it is time dependent in the background. Non-Newtonian density perturbations in \mathcal{G} on superhorizon scales may have an entirely different time dependence in \mathcal{G}' , and the choice of gauge transformation $\mathcal{G} \rightarrow \mathcal{G}'$ is quite arbitrary.

But arbitrariness does not imply that one gauge is correct and all others wrong. Rather, it imposes on physicists the requirement to agree on a convention, otherwise there will be problems in the interpretation of results. The formalism we chose in Chapter 3, which led to the Einstein Equation (3.29) and to Friedmann's Equations (5.4) and (5.5), implicitly used a conventional gauge. Alternatively one could have used gauge-invariant variables, but at the cost of a very heavy mathematical apparatus. Another example concerns the electroweak theory, in which particle states are represented by gauge fields that are locally gauged.

10.2 Structure Formation

As we have seen in the FLRW model, the force of gravity makes a homogeneous matter distribution unstable: it either expands or contracts. This is true for matter on all scales, whether we are considering the whole Universe or a tiny localized region. But the FLRW expansion of the Universe as a whole is not exponential and therefore it is too slow to produce our Universe in the available time. This requires a different mechanism to give the necessary exponential growth: cosmic inflation. Only after the graceful exit from inflation does the Universe enter the regime of Friedmann expansion, during which the Hubble radius gradually overtakes the inflated regions.

In Figure 10.1 we show schematically the fate of fluctuations as a function of time or the cosmic scale a . Inflationary fluctuations crossed the post-inflationary Hubble radius at scale a_1 and came back into vision recently, at $\approx a_0$, with a wavelength λ_{hor} corresponding to the size the Hubble radius at that moment. Some galaxies will cross the post-inflationary Hubble radius at scale a_2 and come back into vision with a wavelength λ_{gal} later at the cosmic scale $> a_H$ with a comoving length scale $> H^{-1}$.

In noninflationary cosmology, a given scale crosses the horizon but once, while in the inflationary cosmology all scales begin subhorizon sized, cross outside the Hubble radius during inflation, and re-enter during the post-inflationary epoch. The largest scales cross outside the Hubble radius first and re-enter last. Causal microphysics operates only on scales less than H^{-1} below the strong black line. During inflation H^{-1} is a constant, and in the post-inflation era it is proportional to $a^{1/n}$, where $n = 2$ during radiation domination, and $n = \frac{3}{2}$ during matter domination.

Jeans Mass. Primordial density fluctuations expand linearly at a rate slower than the Universe is expanding in the mean, until eventually they reach a maximum size and collapse nonlinearly. If the density fluctuates locally, also the cosmic scale factor will be a fluctuating function $a(\mathbf{r}, t)$ of position and time. In overdense regions where the gravitational forces dominate over pressure forces, causing matter to contract

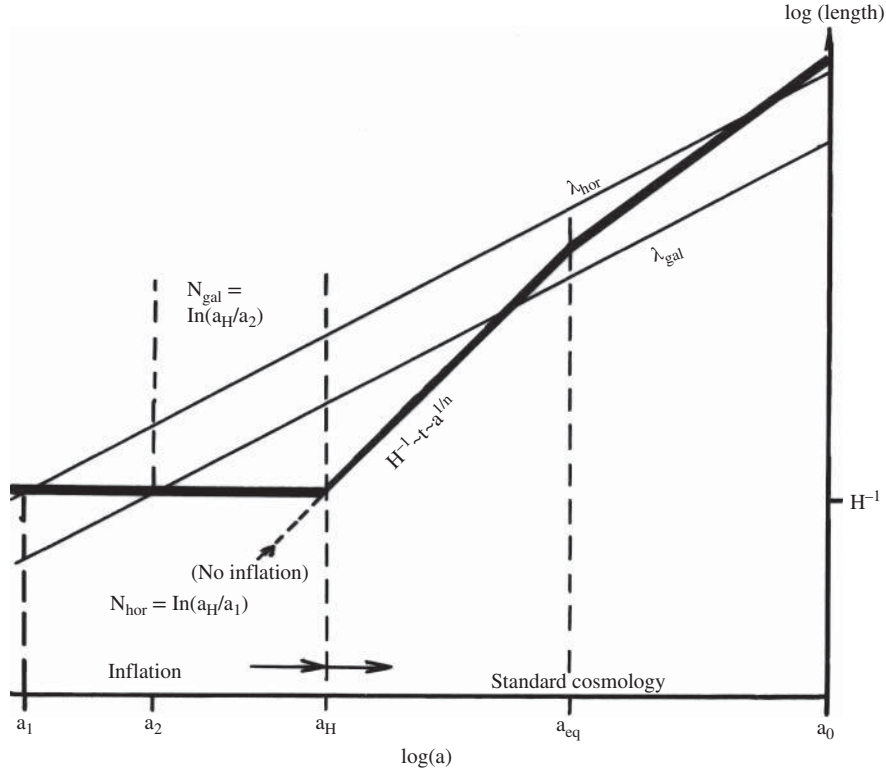


Figure 10.1 The comoving length scale of fluctuations versus the cosmic scale of the Universe a . The wavelength λ_{hor} corresponds to fluctuations which disappeared over the horizon at a_1 and re-entered only recently, at $\approx a_0$. The growth in the scale factor a while the fluctuations are outside the horizon, from a_1 to the end of inflation at a_H , is for our present horizon, $N_{\text{hor}} = \ln(a_H/a_1) \approx 53$. Fluctuations exiting at a_2 will correspond to galaxies entering at a_H , when the scale factor has grown by $N_{\text{gal}} = \ln(a_H/a_2) \approx 45$.

locally and to attract surrounding matter which can be seen as inflow. In other regions where the pressure forces dominate, the fluctuations move as sound waves in the fluid, transporting energy from one region of space to another.

The dividing line between these two possibilities can be found by a classical argument. Let the time of free fall to the center of an inhomogeneity in a gravitational field of strength G be

$$t_G = 1/\sqrt{G\rho}. \tag{10.19}$$

Sound waves in a medium of density ρ and pressure p propagate with velocity

$$c_s = \sqrt{\frac{\partial p}{\partial \rho}},$$

so they move one wavelength λ in the time

$$t_s = \lambda/c_s. \tag{10.20}$$

Note that only baryonic matter experiences pressure forces. Recall that dark matter is pressureless, feeling only gravitational forces.

When t_G is shorter than t_s , the fluctuations are unstable and their amplitude will grow by attracting surrounding matter, becoming increasingly unstable until the matter eventually collapses into a gravitationally bound object. The opposite case is stable: the fluctuations will move with constant amplitude as sound waves. Setting $t_G = t_s$, we find the limiting *Jeans wavelength* $\lambda = \lambda_J$ at the *Jeans instability*, discovered by Sir *James Jeans* (1877–1946) in 1902,

$$\lambda_J = \sqrt{\frac{\pi}{G\rho}} c_s. \quad (10.21)$$

Actually, the factor $\sqrt{\pi}$ was not present in the above Newtonian derivation; it comes from an exact treatment of Equations (10.1)–(10.3); see for instance [3]. The mass contained in a sphere of radius λ_J is called the *Jeans mass*,

$$M_J = \frac{4}{3}\pi\lambda_J^3\rho. \quad (10.22)$$

In order for tiny density fluctuations to be able to grow to galactic size, there must be enough time, or the expansion must be exponential in a brief time. The gravitational collapse of a cloud exceeding the Jeans mass develops exponentially, so the cloud halves its radius in equal successive time intervals. But galaxies and large-scale structures do not condense out of the primordial medium by exponential collapse. The structures grow only linearly with the scale a or as some low power of a .

For subhorizon modes, the distinction between the radiation-dominated and matter-dominated eras is critical. During the radiation era, growth of perturbations is suppressed. During the matter era, perturbations can grow. But during the matter era the Jeans wavelength provides an important boundary. Large wavelength fluctuations will grow with the expansion as long as they are in the linear regime. In an accelerated expansion driven by dark energy, the condition for gravitational collapse becomes extremely complicated. This happens rather late, only when matter domination ends and dark energy becomes dynamically important ($z \sim 1$).

For wavelengths less than the Jeans wavelength the pressure in the baryonic matter can oppose the gravitational collapse and perturbations will oscillate in the linear regime as sound waves, never reaching gravitational collapse. An alternative way of stating this is to note that the radiation pressure and the tight coupling of photons, protons and electrons causes the fluid to be viscous. On small scales, photon diffusion and thermal conductivity inhibit the growth of perturbations as soon as they arise, and on large scales there is no coherent energy transport.

Mass fluctuations at still shorter wavelength, with $\lambda \approx r_U \ll r_H$, can break away from the general expansion and collapse to bound systems of the size of galaxies or clusters of galaxies. Fluctuations which enter in the nonlinear regime, where the ratio in Equation (10.17) is large, collapse rapidly into black holes before pressure forces have time to respond.

For baryonic matter before the recombination era, the baryonic Jeans mass is some 30 times larger than the mass M_H of baryons within the Hubble radius r_H , so if there

exist nonlinear modes they are outside it (the Jeans wavelength is greater than the horizon). A mass scale M is said to enter the Hubble radius when $M = M_H$. Well inside the Hubble radius, the fluctuations may start to grow as soon as the Universe becomes matter dominated, which occurs at time $t_{\text{eq}} \approx 54\,500$ yr.

Upon recombination, the baryonic Jeans mass falls dramatically. If the fluid is composed of some nonbaryonic particle species (cold dark matter), the Jeans wavelength is small after radiation–matter equality, allowing subhorizon perturbations to grow from this time. After matter–radiation equality, nonbaryonic matter can form potential wells into which baryons can fall after recombination.

Matter can have two other effects on perturbations. Adiabatic fluctuations lead to gravitational collapse if the mass scale is so large that the radiation does not have time to diffuse out of one Jeans wavelength within the time t_{eq} . As the Universe approaches decoupling, the photon mean free path increases and radiation can diffuse from overdense regions to underdense ones, thereby smoothing out any inhomogeneities in the plasma. For wavelengths below the Jeans wavelength, *collisional dissipation* or *Silk damping* (after *J. Silk*) erases perturbations in the matter (baryon) radiation field through photon diffusion. This becomes most important around the time of recombination. Random waves moving through the medium with the speed of sound c_s erase all perturbations with wavelengths less than $c_s t_{\text{eq}}$. This mechanism sets a lower limit to the size of the structures that can form by the time of recombination: they are not smaller than rich clusters or superclusters. But, in the presence of nonbaryonic matter, Silk damping is of limited importance because nonbaryonic matter does not couple with the radiation field.

The second effect is free streaming of weakly interacting relativistic particles such as neutrinos. This erases perturbations up to the scale of the horizon, but this also ceases to be important at the time of matter–radiation equality.

The situation changes dramatically at recombination, when all the free electrons suddenly disappear, captured into atomic Bohr orbits, and the radiation pressure almost vanishes. This occurs at time 400 000 yr after Big Bang (see Figure 6.5). Now the density perturbations which have entered the Hubble radius can grow with full vigour.

Sunyaev–Zel’dovich Effect (SZE). At some stage the hydrogen gas in gravitationally contracting clouds heats up enough to become ionized and to re-ionize the CMB: the Sunyaev–Zel’dovich effect. We refer to the Planck result in Equation (8.46) that such re-ionization clouds occur at a half-reionization redshift $z_r \approx 11.15$.

The free electrons and photons in the ionized clouds build up a radiation pressure, halting further collapse. The state of such clouds today depends on how much mass and time there was available for their formation. Small clouds may shrink rapidly, radiating their gravitational binding energy and fragmenting. Large clouds shrink slowly and cool by the mechanism of electron Thomson scattering. As the recombination temperature is approached the photon mean free paths become larger, so that radiation can diffuse out of overdense regions. This damps the growth of inhomogeneities.

The distortion of the CMB spectrum due to the SZE can be used to detect intergalactic clouds and to provide another estimate of H_0 by combining radio and X-ray observations to obtain the distance to the cloud. The importance of the SZE surveys is that they are able to detect all clusters above a certain mass limit independent of the redshifts of the clusters. The ratio of the magnitude of the SZE to the CMB does not change with redshift. The effects of re-ionization on the CMB temperature–polarization power were discussed in Section 8.4.

Structure Sizes and Formation Times. Only clouds exceeding the Jeans mass stabilize and finally attain *virial equilibrium*. It is intriguing (but perhaps an accident) that the Jeans mass just after recombination is about $10^5 M_\odot$, the size of globular clusters! Galaxies have masses of the order of $10^{12} M_\odot$ corresponding to fluctuations of order $\delta \simeq 10^{-4}$ as they cross the horizon. We have already made use of this fact to fix the mass m_ϕ of the scalar field in Equation (7.44).

The timetable for galaxy and cluster formation is restricted by two important constraints. At the very earliest, the Universe has to be large enough to have space for the first formed structures. If these were galaxies of the present size, their number density can be used to estimate how early they could have been formed. We leave this for a problem.

The present density of the structures also sets a limit on the formation time. The density contrast at formation must have exceeded the mean density at that time, and since then the contrast has increased with a^3 . Thus, rich clusters, for instance, cannot have been formed much earlier than at

$$1 + z \approx 2.5 \Omega^{-1/3}. \quad (10.23)$$

It seems that all the present structure was already in place at $z = 5$. This does not exclude the fact that the largest clusters are still collapsing today. In a critical universe structure formation occurs continuously, rich galaxy clusters form only at a redshift of 0.2–0.3, and continue to accrete material even at the present epoch. In that case many clusters are expected to show evidence for recent merger events and to have irregular morphologies. There is clear observational evidence that star formation activity drives gas out of galaxies.

As a result of mass overdensities, the galaxies influenced by the ensuing fluctuations in the gravitational field will acquire peculiar velocities. One can derive a relation between the mass autocorrelation function and the RMS peculiar velocity (see reference [3]). If one takes the density contrast to be $\delta_m = 0.3$ for RMS fluctuations of galaxy number density within a spherical volume radius $30 h^{-1}$ Mpc, and if one further assumes that all mass fluctuations are uncorrelated at larger separations, then the acceleration caused by the gravity force of the mass fluctuations would predict deviations from a homogeneous peculiar velocity field in rough agreement with observations in our neighborhood. Much larger density contrast would be in marked disagreement with the standard model and with the velocity field observations.

Problems

1. The mean free path ℓ of photons in homogeneous interstellar dust can be found from Equation (1.4) assuming that the radius of dust grains is 10^{-7} m. Extinction observations indicate that $\ell \approx 1$ kpc at the position of the Solar System in the Galaxy. What is the number density of dust grains [4]?
2. To derive Jeans wavelength λ_J and Jeans mass M_J [see Equation (10.22)], let us argue as follows. A mass $M_J = \rho \lambda_J^3$ composed of a classical perfect gas will collapse gravitationally if its internal pressure $P = \rho kT/m$ cannot withstand the weight of a column of material of unit area and height λ_J . Here m is the mean mass of the particles forming the gas. If we set the weight of the latter greater than or equal to P ,

$$\frac{GM_J}{\lambda_J^2} \rho \lambda_J \gtrsim \frac{\rho kT}{m},$$

we will obtain a constraint on the sizes of fragment which will separate gravitationally out of the general medium. Show that this leads to Equation (10.22) [5].

3. The universal luminosity density radiated in the blue waveband by galaxies is

$$L_U = (2 \pm 0.2) \times 10^8 h L_\odot \text{ Mpc}^{-3}.$$

Show that the Coma value $Y = M/L = 300$ in solar units then gives $\Omega_m = 0.30$.

4. Assuming that galaxies form as soon as there is space for them, and that their mean radius is $30 h^{-1}$ kpc and their present mean number density is $0.03 h^3 \text{ Mpc}^{-3}$, estimate the redshift at the time of their formation [2].

References

- [1] Planck Collaboration: Ade, P. A. R. *et al.* 2014 *Astron. Astrophys* and preprint arXiv: 1303.5076 [astro-ph.CO].
- [2] Hawkins, E. *et al.* 2003 *Mon. Not. R. Astron. Soc.* **346**, 78.
- [3] Peebles, P. J. E. 1993 *Principles of physical cosmology*. Princeton University Press, Princeton, NJ.
- [4] Shu, F. H. 1982 *The physical universe*. University Science Books, Mill Valley, CA.
- [5] York, D. G. *et al.* 2000 *Astr. J.* **120**, 1579.

Dark Energy

The *Friedmann–Lemaître universe* which we met in Chapter 5 has become the generally accepted description of the Universe, the *FLRW concordance model*. In Chapter 8 we learned the reason for this: all observations show that the Universe is expanding and characterized by a positive *cosmological constant* λ and a large density parameter, Ω_λ ; see Equations (5.20), (8.50). This is also clear from Figure 8.7.

What λ stands for is not known. Dark energy came as a complete surprise. Nothing in big bang or inflationary cosmology predicted its existence. Is the Universe dominated by some new form of dark energy or does Einstein's theory of gravity break down on cosmological scales? In the first case λ belongs to the energy-momentum tensor on the right side of the Einstein Equation (3.29), in the second case it is a modification of the geometry described by the left side of Equation (3.29). The introduction in Section 11.2 of a new single field could be a modification of the geometry or a fluid added to the energy-momentum tensor. We spend most time on quintessence here.

In Section 11.3 we introduce $f(R)$ models which are modifications of the geometry of spacetime. Even more exotic modifications of the geometry of spacetime are models defined in extra dimensions in Section 11.4. We also list a number of other models here.

11.1 The Cosmological Constant

The introduction of the cosmological constant into our description of the Universe is problematic for at least three reasons. First, as we noted in Equation (5.22), its present value is extremely small, in fact some 122 orders of magnitude smaller than theoretical expectations. The density is about

$$\rho_\lambda \approx 2.9 \times 10^{-47} \text{ GeV}^4.$$

If ρ_λ were even slightly larger, the repulsive force would cause the Universe to expand too fast so that there would not be enough time for the formation of galaxies or other gravitationally bound systems. This is called the *cosmological constant problem*.

Second, it raises the question of why the sum

$$\Omega_0 = \Omega_m + \Omega_\lambda$$

is precisely 1.0 today (to within 0.3%) when we are there to observe it, after an expansion of some 12 billion years when it was always greater than 1.0. The density of matter decreases like a^{-3} , while Ω_λ remains constant, so why has the cosmological constant been fine-tuned to come to dominate the sum only now? This is referred to as the *cosmic coincidence problem*.

Third, we do not have the slightest idea what the λ energy consists of, only that it distorts the geometry of the Universe as if it were matter with strongly negative pressure, and it acts as an *anti-gravitational* force which is unclustered at all scales. Note that a positive λ in Equation (5.18) implies an accelerated expansion, $\ddot{a} > 0$. Since we know so little about it, we also cannot be sure that λ is constant in time, and that its equation of state is always $w_\lambda = -1$.

Observations. The first evidence for the need of a cosmological constant came in 1998–1999 from two independent teams monitoring high-redshift, Type Ia supernovae [1, 2]. The expectation was to see cosmic deceleration, since gravitation is attractive. But the teams converged on the remarkable result that, on the contrary, the cosmic expansion was accelerating, consistent with a flat universe with $\Omega_\lambda \approx 0.7$. When compared to local Type Ia supernovae, those observed at $z \approx 0.5$ were fainter than expected in a matter-dominated universe, as if their light were absorbed by intervening dust. However, many systematic checks ruled out the hypothesis of grey dust extinction that increased towards higher redshifts. While the significance of this discovery has since then only been confirmed with the inclusion of larger and better calibrated SN Ia data sets the cause of the acceleration remains unknown.

At the time of the first discovery by the supernova teams the ground had been well prepared by the CMB and large scale structure data, which already provided substantial indirect evidence for a cosmological constant. As is evident from the very recent Figure 8.7, the constraints from the supernovae and from CMB are practically orthogonal in the Ω_m, Ω_λ space.

The supernova results were followed within a year by the results of balloon-borne CMB experiments that mapped the first acoustic peak and measured its angular location, providing strong evidence for spatial flatness [3, 4]. The acoustic peak measurement implied that the alternative to λ was not an open universe but a strongly decelerating, $\Omega_m = 1$ universe that disagreed with the supernova data by 0.5 magnitudes, a level much harder to explain with observational or astrophysical effects. The combination of spatial flatness and improving measurements of the Hubble constant provided an entirely independent argument for an energetically dominant accelerating component: a matter-dominated universe with $\Omega_{tot} = 1$ would have age $t_0 = (2/3)H_0 \approx 9.5$ Gyr, too young to accommodate the 12–14 Gyr ages estimated for globular clusters.

Explaining simultaneously all these and still other data, all consistent with an inflationary cold dark matter model with a cosmological constant, requires an accelerating universe. Recent supernova data plotted in Figure 11.1 as $H(z)/(1+z)$ versus z [5] show that an early deceleration followed by a recent acceleration is favored, determining the time of transition from deceleration to acceleration to $z = 0.74 \pm 0.05$. Observations of baryonic acoustic oscillations [6] plotted in Figure 11.2 as $H(z)/(1+z)$ versus z determine this moment to have occurred at $z \approx 0.8$.

Dynamical Models. The cosmological constant corresponds to static gravity with the equation of state $w = -1$ as in Equation (5.28). If w were a function of the scale a the expansion history would be different, but unfortunately all functions are *ad hoc*. Some simple one-parameter formulas are

$$w(z) = w_0 + w_1 z, \tag{11.1}$$

and

$$w(a) = w_0 + w_1(1 - a). \tag{11.2}$$

One relatively successful two-parameter formula is

$$w(a) = w_0[1 + b \ln(1 + z)]^{-2}. \tag{11.3}$$

The advantage of this formula is that it covers the whole observed range of z .

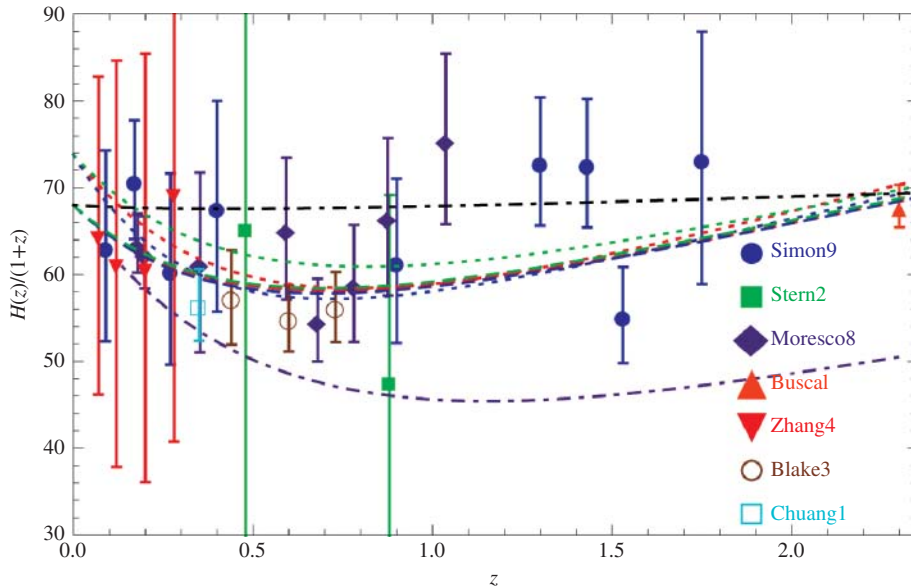


Figure 11.1 Evidence for transition from deceleration in the past to acceleration today. From reference [6]. From Farook, O. and Ratra, B., Hubble parameter measurement constraints on the cosmological deceleration-acceleration transition redshift, *Astrophys. J. Lett.*, **766**, L7, published 4 March 2013. © AAS. Reproduced with permission. (See plate section for color version.)

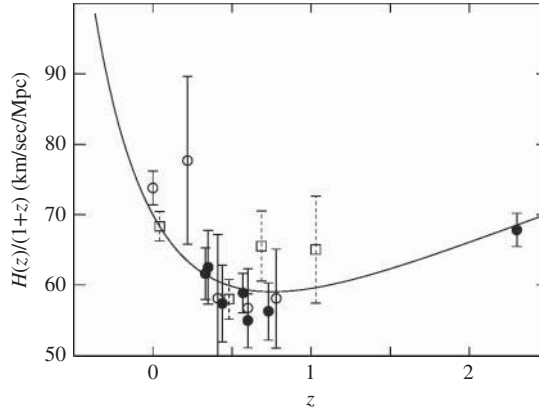


Figure 11.2 Measurements of $H(z)/(1+z)$ versus z demonstrating the acceleration of the expansion for $z < 0.8$ and deceleration for $z > 0.8$. The BAO-based measurements from different groups are the filled circles. The line is the Λ CDM prediction for $(h, \Omega_m, \Omega_\lambda = 0.7, 0.27, 0.73)$. For further details, see Busca *et al.* [6]. Reproduced with permission © ESO.

Decaying Cosmological Constant. A dynamical approach to remove or alleviate the extreme need for fine-tuning λ is to choose it to be a slowly varying function of time, $\lambda(t)$. The initial conditions require $\lambda(t_{\text{Planck}}) \approx 10^{122} \lambda_0$, from which it decays to its present value at time t_0 .

The Universe is then treated as a fluid composed of dust and dark energy in which the dark energy density, $\rho_\lambda(t) = \lambda(t)/8\pi G$, continuously transfers energy to the material component. Its equation of state is then of the form

$$p_\lambda = -\rho_\lambda \left[1 + \frac{1}{3} \frac{d \ln \rho_\lambda(a)}{d \ln a} \right]. \quad (11.4)$$

In the classical limit when $\rho_\lambda(a)$ is a very slow function of a so that the derivative term can be ignored, one obtains the equation of state of the cosmological constant, $w_\lambda = -1$.

The advantage in removing the need for fine-tuning is, however, only replaced by another arbitrariness: an ansatz for $\lambda(t)$ is required and new parameters characterizing the timescale of the deflationary period and the transfer of energy from dark energy to dust must be introduced. Such phenomenological models have been presented in the literature [7, 8], and they can lead to testable predictions.

11.2 Single Field Models

Instead of arguing about whether λ should be interpreted as a correction to the geometry or to the stress–energy tensor, we could go the whole way and postulate the existence of a new kind of energy, described by a slowly evolving scalar field $\varphi(t)$ that contributes to the total energy density together with the background (matter and radiation) energy density. This scalar field is assumed to interact only with gravity and

with itself. In Section 7.2 we have already discussed this situation in the context of single-field inflation.

Since a scalar field is mathematically equivalent to a fluid with a time-dependent speed of sound, one can find potentials $V(\varphi)$ for which the dynamical vacuum acts like a fluid with negative pressure, and with an energy density behaving like a decaying cosmological constant. In comparison with plain $\lambda(t)$ models, scalar field cosmologies have one extra degree of freedom, since both a potential $V(\varphi)$ and a kinetic term $\frac{1}{2}\dot{\varphi}^2$ need to be determined.

The simplest example of a minimally coupled and spatially homogeneous scalar field has the Lagrangian density 5.86 and the equation of motion is then given by the *Klein–Gordon* Equation 7.24, where the prime indicates derivation with respect to φ .

The energy density and pressure enter in the diagonal elements of $T_{\mu\nu}$, and they are

$$\rho_\varphi c^2 = \frac{1}{2}\dot{\varphi}^2 + V(\varphi) \quad \text{and} \quad p_\varphi = \frac{1}{2}\dot{\varphi}^2 - V(\varphi), \quad (11.5)$$

respectively. Clearly the pressure is always negative if the evolution is so slow that the kinetic energy density $\frac{1}{2}\dot{\varphi}^2$ is less than the potential energy density. Note that in Equations (7.24) and (11.5) we have ignored terms describing spatial inhomogeneity which could also have been present.

The conservation of energy-momentum for the scalar field is as in Equation (5.24),

$$\dot{\rho}_\varphi + 3H\rho_\varphi(1 + w_\varphi) = 0. \quad (11.6)$$

As in Equation (5.29), the energy density of the scalar field decreases as $a^{-3(1+w_\varphi)}$. Inserting Equations (11.5) into Equation (11.6), one indeed obtains Equation (7.24). The equation of state of the φ field is then a function of the cosmological scale a (or time t or redshift z),

$$w_\varphi = \frac{\dot{\varphi}^2 + 2V(\varphi)}{\dot{\varphi}^2 - 2V(\varphi)}. \quad (11.7)$$

Starting early from a wide range of arbitrary initial conditions w_φ oscillates between ≈ 1 and ≈ -1 , until an epoch when ρ_φ freezes to a small value.

The conditions for acceleration are

$$w_\varphi < -1/3, \quad a(t) \propto t^d \quad \text{with } d > 1$$

so that

$$p_\varphi < 0 \quad \text{or} \quad \rho_\varphi \propto a^{-2}.$$

However, dark energy defined this way and called *quintessence* turns out to be another *Deus ex machina* which not only depends on the parametrization of an arbitrary function $V(\varphi)$, but also has to be fine-tuned initially in a way similar to the cosmological constant.

The Inflaton as Quintessence. Now we have met two cases of scalar fields causing expansion: the inflaton field acting before t_{GUT} and the quintessence field describing present-day dark energy. It would seem economical if one and the same scalar field

could do both jobs. Then the inflaton field and quintessence would have to be matched at some time later than t_{GUT} . This seems quite feasible since, on the one hand, the initially dominating inflaton potential $V(\varphi)$ must give way to the background energy density $\rho_r + \rho_m$ as the Universe cools, and on the other hand, the dark energy density must have been much smaller than the background energy density until recently. Recall that quintessence models are constructed to be quite insensitive to the initial conditions.

On the other hand, nothing forces the identification of the inflaton and quintessence fields. The inflationary paradigm in no way needs nor predicts quintessence.

In the previously described models of inflation, the inflaton field φ settled to oscillate around the minimum $V(\varphi = 0)$ at the end of inflation. Now we want the inflaton energy density to continue a monotonic roll-down toward zero, turning ultimately into a minute but nonvanishing quintessence tail. The global minimum of the potential is only reached in a distant future, $V(\varphi \rightarrow \infty) \rightarrow 0$. In this process the inflaton does not decay into a thermal bath of ordinary matter and radiation because it does not interact with particles at all, it is said to be sterile. A sterile inflaton field avoids violation of the equivalence principle, otherwise the interaction of the ultralight quintessence field would correspond to a new long-range force. Entropy in the matter fields comes from gravitational generation at the end of inflation rather than from decay of the inflaton field.

The task is then to find a potential $V(\varphi)$ such that it has two phases of accelerated expansion: from t_{P} to t_{end} at the end of inflation, and from a time $t_{\text{F}} \approx t_{\text{GUT}}$ when the instanton field freezes to a constant value until now, t_0 . Moreover, the inflaton energy density must decrease faster than the background energy density, equalling it at some time t_* when the field is φ_* , and thereafter remaining subdominant to the energy density of the particles produced at t_{end} . Finally it must catch up with a tracking potential at some time during matter domination, $t > t_{\text{eq}}$.

The mathematical form of candidate potentials is of course very complicated, and it would not be very useful to give many examples here. However, it is instructive to follow through the physics requirements on φ and $V(\varphi)$ from inflation to present.

Kination. Inflation is caused by an essentially constant potential $V(\varphi)$ according to Equation (7.36). The condition $V(\varphi \rightarrow \infty) \rightarrow 0$ requires an end to inflation at some finite time t_{end} when the field is φ_{end} and the potential is $V_{\text{end}} \equiv V(\varphi_{\text{end}})$. The change in the potential at t_{end} from a constant to a decreasing roll then implies, by Equation (7.37), that $\dot{\varphi}_{\text{end}} \neq 0$, and furthermore, by Equation (7.24), that also $\ddot{\varphi}_{\text{end}} \neq 0$. Then the slow-roll conditions in Equation (7.32) for ϵ and η are also violated.

During inflation the kinetic energy density of the inflaton is

$$\rho_{\text{kin}} = \epsilon V = \frac{m_{\text{Planck}}^2}{16\pi} \left[\frac{V'^2(\varphi)}{V(\varphi)} \right]. \quad (11.8)$$

Thus when $V'(\varphi)$ starts to grow, so does ρ_{kin} , and the total energy density of the Universe becomes dominated by the inflaton kinetic energy density. This epoch has been

called *kination* or *deflation*. Equation (11.7) then dictates that the equation of state is

$$w_\varphi = \frac{\dot{\varphi}^2 + 2V(\varphi)}{\dot{\varphi}^2 - 2V(\varphi)} \approx 1, \quad (11.9)$$

so that the kinetic energy density decreases as

$$\rho(a) \propto a^{-3(1+w)} = a^{-6} \quad (11.10)$$

from Equation (5.29). This is much faster than the a^{-4} decrease of the radiation energy density ρ_r , and the a^{-3} decrease of the initially much smaller matter energy density ρ_m . Consequently, kination ends at the moment when ρ_r overtakes ρ_{kin} at time t_* . When constructing phenomenological models for this scenario, one constraint is of course that $\rho_r(t_{\text{end}}) \ll V_{\text{end}}$, or equivalently, $t_{\text{end}} < t_*$. This behavior is well illustrated in Figure 11.3, taken from the work of Dimopoulos and Valle [9].

Since matter and radiation are gravitationally generated at t_{end} , the reheating temperature of radiation is given by

$$T_{\text{reh}} = \alpha T_{\text{H}}, \quad (11.11)$$

where T_{H} is the Hawking temperature [Equation (3.35)], and α is some reheating efficiency factor less than unity. In Figure 11.3 the radiation energy density $\rho_r \equiv \rho_r$ starts at $T_{\text{reh}}^4 \ll V_{\text{end}}$, and then catches up with $V(\varphi)$ at φ_* . Now the Universe becomes radiation dominated and the hot Big Bang commences. Note that the term ‘hot Big Bang’ has a different meaning here: it does not refer to a time zero with infinite temperature, but to a moment of explosive entropy generation. This mimics the Big Bang so that all of its associated successful predictions ensue.

Quintessence. The properties of the inflaton field plays no role any more during the Big Bang, they are fixed by requirements at Planck time when the quintessence field is completely negligible. The properties of quintessence, on the other hand, are fixed by present observations. This makes their identification rather artificial.

Let us continue the argument of the previous paragraph. The kinetic energy density reduces rapidly to negligible values by its a^{-6} dependence and the field freezes ultimately to a nonzero value φ_{F} . The residual inflaton potential $V(\varphi)$ again starts to dominate over the kinetic energy density, however, staying far below the radiation energy density and, after t_{eq} , also below the matter energy density.

As we approach t_0 , the task of phenomenology is to devise a quintessence potential having a suitable tracker. The nature of the tracker potential is decided by the form of the quintessence potential. To arrive at the present-day dark energy which causes the evolution to accelerate, the field φ must be unfrozen again, so φ_{F} should not be very different from φ_0 . Many studies have concluded that only exponential trackers are admissible, and that quintessence potentials can be constructed by functions which behave as exponentials in φ early on, but which behave more like inverse power potentials in the quintessential tail. A simple example of such a potential is

$$V(\varphi \gg \varphi_{\text{end}}) \approx V_{\text{end}} \frac{\exp(-\lambda\varphi/m_{\text{P}})}{(\varphi/m)^k}, \quad (11.12)$$

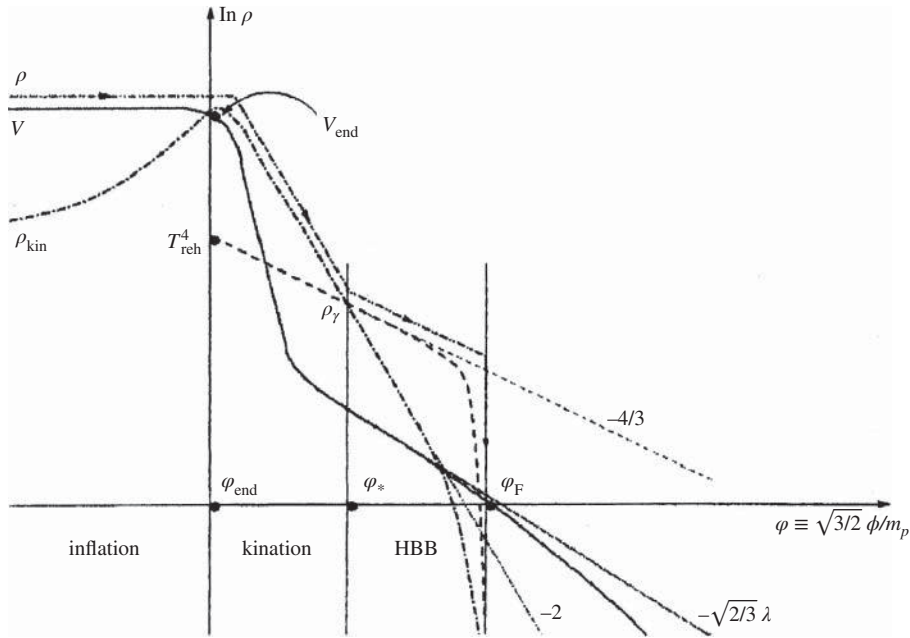


Figure 11.3 Schematic view of the scalar potential from inflation to quintessence. The potential V (solid line) features two flat regions, the inflationary plateau and the quintessential tail. The inflation is terminated at φ_{end} by a drastic reduction of V , leading to a rapid roll-down of the scalar field from the inflationary plateau towards the quintessential tail. At the end of inflation the kinetic energy density of the scalar field, ρ_{kin} (dash-dotted line), dominates for a brief period the energy density of the Universe. During this time the radiation energy density ρ_γ (dashed line) reduces less rapidly and catches up with ρ_{kin} at time t_* when the field is φ_* , and the explosive generation of entropy commences. After that the kinetic energy of the scalar field reduces rapidly to zero and the field freezes asymptotically to a value φ_F , while the overall energy density of the Universe (dash-dot-dotted line) continues to decrease due to the Hubble expansion. Assuming a quasi-exponential tail given by Equation (11.12), the potential beyond φ_F is seen departing logarithmically from a pure exponential case (dotted line). Reprinted from K. Dimopoulos and J. W. F. Valle, *Modeling quintessential inflation* [9], copyright 2002, with permission from Elsevier.

where $k \geq 1$ is an integer, $\lambda > 0$ is a parameter characterizing the exponential scale, and $m < m_p$ is a mass scale characteristic of the adopted inverse power-law scale. For more details, see Dimopoulos and Valle [9].

Tracking Quintessence. In a somewhat less arbitrary model [10, 11], one constructs quintessence in such a way that its energy density is smaller than the background component for most of the history of the Universe, somehow tracking it with the same time dependence. This field converges onto the tracker solution with nearly constant slope $-d \ln \rho_\phi / d \ln a$ at a value $< d \ln \rho_m / d \ln a$, and regardless of the value of w_ϕ , in the radiation-domination epoch, w_ϕ automatically decreases to a negative value at time t_{eq} when the Universe transforms from radiation domination to matter domination.

We saw in Equation (5.33) that radiation energy density evolves as a^{-4} —faster than matter energy density, a^{-3} . Consequently, ρ_r is now much smaller than ρ_m .

But once w_ϕ is negative, ρ_ϕ decreases at a slower rate than ρ_m so that it eventually overtakes it. At that moment, $\dot{\phi}(t)$ slows to a near stop, causing w_ϕ to decrease toward -1 , and tracking stops. Judging from the observed large value of the cosmological constant density parameter today, $\Omega_\Lambda = 0.714$, this happened in the recent past when the redshift was $z \sim 2\text{--}4$. Quintessence is already dominating the total energy density, driving the Universe into a period of de Sitter-like accelerated expansion.

The tracker field should be an attractor in the sense that a very wide range of initial conditions for ϕ and $\dot{\phi}$ rapidly approach a common evolutionary track, so that the cosmology is insensitive to the initial conditions. Such cosmological solutions, called *scaling solutions*, satisfy $\rho_\phi(t)/\rho_m(t) = \text{constant} > 0$. Scaling solutions define the borderline between deceleration and acceleration.

Thus the need for fine-tuning is entirely removed, the only arbitrariness remains in the choice of a function $V(\phi)$. With a judicious choice of parameters, the coincidence problem can also be considered solved, albeit by tuning the parameters ad hoc.

In Chapter 7 we already discussed the inflationary de Sitter expansion following the Big Bang, which may also be caused by a scalar *inflaton field*. Here we just note that the initial conditions for the quintessence field can be chosen, if one so desires, to match the inflaton field.

Tracking behavior with $w_\phi < w_b$ occurs [10, 11] for any potential obeying

$$\Gamma \equiv V''V/(V')^2 > 1, \quad (11.13)$$

and which is nearly constant over the range of plausible initial ϕ ,

$$\frac{d(\Gamma - 1)}{H dt} \ll |\Gamma - 1|, \quad (11.14)$$

or if $-V'/V$ is a slowly decreasing function of ϕ . Many potentials satisfy these criteria, for instance power law, exponential times power law, hyperbolic, and Jacobian elliptic functions. For a potential of the generic form,

$$V(\phi) = V_0(\phi_0/\phi)^{-\beta} \quad (11.15)$$

with β constant, one has a good example of a tracker field for which the kinetic and potential terms remain in a constant proportion.

The values of w_ϕ and Ω_ϕ depend both on $V(\phi)$ and on the background. The effect of the background is through the $3H\dot{\phi}$ term in the scalar field equation of motion [Equation (7.24)] when w changes, H also changes, which, in turn, changes the rate at which the tracker field evolves down the potential.

The tracking potential is characterized as *slow rolling* when the slow-roll parameters [already defined in Equation (7.34)]

$$\eta(\phi) \equiv \frac{m_{\text{Planck}}^2}{16\pi} \left(\frac{V''}{V} \right) \ll 1, \quad \epsilon \equiv \frac{m_{\text{Planck}}^2}{16\pi} \left(\frac{V'}{V} \right)^2 \ll 1, \quad (11.16)$$

meaning that $\ddot{\phi}$ in Equation (7.24) and $\dot{\phi}^2$ in Equation (11.5) are both negligible. At very early times, however, $-V'/V$ is slowly changing, but is itself not small. This establishes the important distinction between static and quasi-static quintessence with

$w_\varphi \approx -1$ and dynamical quintessence with $w_\varphi > -1$. This means that the slow-roll approximation is not necessarily applicable to dynamical quintessence, and that the latter generally requires exact solution of the equation of motion [Equation (7.24)].

Given a potential like Equation (11.15) and fixing the current values of parameters $\Omega_m, \Omega_r, \Omega_\varphi, w_\varphi$ one can solve the equation of motion [Equation (7.24)] by numerical integration. Finding the functions $\Gamma(a), w_\varphi(\varphi), w_\varphi(a)$ or $\Omega_\varphi(a)$ is a rather complicated exercise. Finally, when the potential gets shallow, the relatively fast-rolling dynamical quintessence solution exits from the scaling regime and becomes static, approaching the $w_\varphi = -1$ regime, and an accelerated expansion commences.

The equation of state of dark energy, w_φ , introduces a new degeneracy with Ω_m and h which cannot be resolved by CMB data alone. Using the full set of available data in Figure 9.1, one can obtain limits to w_φ .

In this model the lookback time is given by

$$t(z) = \frac{1}{H_0} \int_1^{1/(1+z)} da [(1 - \Omega_0) + \Omega_m a^{-1} + \Omega_r a^{-2} + \Omega_\varphi a^{2-3(1+w_\varphi)}]^{-1/2}. \quad (11.17)$$

For $1/(1+z) = 0$ this gives us the age [Equation (5.55)] of the Universe, t_0 .

For the case of $-1 < w_\varphi$, the effect on galaxy formation is earlier collapse times and more rapid collapse of overdensities. The equation of state and the amplitude of velocity fluctuations σ_8 are correlated in such a way that for constant w_φ , higher (lower) σ_8 produces earlier (later) collapse times, respectively.

K-Essence. As already mentioned, the weakness of the tracking quintessence model is that the energy density for which the pressure becomes negative is set by an adjustable parameter which has to be fine-tuned to explain the cosmic coincidence. Surely one can do better by adding degrees of freedom, for instance by letting $\varphi(t)$ interplay with the decaying cosmological constant $\lambda(t)$, or with the matter field, or with a second real scalar field $\psi(t)$ which dominates at a different time, or by taking $\varphi(t)$ to be complex, or by choosing a double-exponential potential, or by adding a new type of matter or a dissipative pressure to the background energy density. Surely the present acceleration could have been preceded by various periods of deceleration and acceleration. These are only a few of the alternatives that have been proposed.

One interesting alternative called *K-essence* [12] comes at the cost of introducing a nonlinear kinetic energy density functional of the scalar field. In Equation (5.86) we introduced a Lagrangian which we now slightly generalize to read

$$\mathcal{L}_{phi,X} = \frac{1}{2} g^{\mu\nu} \partial_\mu \phi \partial_\nu \phi - V(\phi) \equiv X - V(\varphi). \quad (11.18)$$

X could be replaced by a positive semi-definite kinetic function of $K(X)$.

If φ describes a perfect fluid with pressure p and density ρ_k the energy-momentum tensor is given by $p = \mathcal{L}$ and

$$\rho_k = [2XK'(X) - K(X)] + V(\varphi) \quad (11.19)$$

where a prime denotes a partial derivative with respect to X . Then $w_\varphi \equiv p/\rho$. A necessary condition that the theory be stable is $\rho > 0$ and further that the speed of sound $\partial p/\partial \rho$ be positive.

Solving the equations one finds that the solution is tracking, just like quintessence. It is also quite insensitive to the initial conditions. The K-field tracks the radiation energy density until t_{eq} , when a sharp transition from positive to negative pressure occurs, with $w_\varphi = -1$ as a consequence. The K-essence density ρ_k then drops below ρ_m and, thereafter, in the matter-dominated epoch the K-field does not track the background at all, it just stays constant. Thus the time of K-essence domination and accelerated expansion simply depends on t_{eq} . However, this is another case of fine-tuning: ρ_k must drop precisely to the magnitude of the present-day ρ_λ .

Phantom Fields. Could the scalar field obey an equation of state with negative kinetic energy violating the weak energy condition $\rho + p \geq 0$?. This can occur in many models but one is trying to avoid it at all costs. The reason for this is serious problems with quantum instabilities. Since $w_\varphi < -1$, such a universe evolves within a finite time to a *Big Rip singularity* when the curvature grows toward ∞ .

However, the singularity can be avoided if the potential has a maximum, for instance

$$V(\varphi) = \frac{V_0}{\cosh(\varphi/\varphi_0)}. \tag{11.20}$$

In contrast to ordinary kinematics, the phantom field evolves toward the top of the potential and crosses over on the other side. It then turns around and again returns to the top, executing damped oscillations across the top until it finally settles at $w_\varphi = -1$.

This model is strongly disfavored by the data in comparison with genuine $w_\varphi = -1$ models.

Tachyon Fields. Somewhat related are tachyon fields which move with superluminal velocity. Since special relativity in four-dimensional spacetime forbids this, tachyon models require rather drastic revisions of general relativity or of spacetime. Speculations have also appeared in the literature that the Universe might have undergone shorter periods of this type. It is well to remember that nothing is known about whether the cosmological constant is indeed constant or whether it will remain so, nor about the future behavior of a quintessence field and its equation of state.

In higher-dimensional spacetimes tachyons may move in the bulk between our four-dimensional *brane* and other branes. The potential may be of the form 11.20 so that the field has a ground state at $w_\varphi = \infty$.

In a flat FRW background one has

$$\rho_\varphi = \frac{V(\varphi)}{\sqrt{1 - (\dot{\varphi}^2)}} \quad \text{and} \quad p_\varphi = -\frac{V(\varphi)^2}{\rho_\varphi}. \tag{11.21}$$

From Friedmann's equations one then obtains

$$\frac{\ddot{a}}{a} = \frac{8\pi G}{3} \frac{V(\varphi)}{\sqrt{1 - \dot{\varphi}^2}} (1 - 2\dot{\varphi}^2). \tag{11.22}$$

Hence the accelerated expansion occurs for $\dot{\varphi}^2 < 2/3$. Note that $w_\varphi = \dot{\varphi}^2 - 1$ is in the range $-1 < w_\varphi < 0$.

Chaplygin Gas. Another simple and well-studied model of dark energy introduces into $T_{\mu\nu}$ the density ρ_φ and pressure p_φ of an ideal fluid with a constant potential $V(\varphi) = A > 0$ called *Chaplygin gas* following Chaplygin's historical work in aerodynamics. Here A has the dimensions of energy density squared. Ordinary matter is assumed not to interact with Chaplygin gas, therefore one has separate continuity equations for the energy densities ρ_m and ρ_φ of the same form as in FLRW geometry, respectively,

$$\dot{\rho} + 3H(\rho + p) = 0. \quad (11.23)$$

Pressureless dust with $p = 0$ then evolves as $\rho_m(a) \propto a^{-3}$.

The barotropic equation of state is

$$p_\varphi = -A/\rho_\varphi. \quad (11.24)$$

The continuity Equation (11.23) is then

$$\dot{\rho}_\varphi + 3H(\rho_\varphi - A/\rho_\varphi) = 0,$$

which integrates to

$$\rho_\varphi(a) = \sqrt{A + B/a^6}, \quad (11.25)$$

where B is an integration constant. Thus this model has two free parameters.

Obviously the limiting behavior of the energy density is

$$\rho_\varphi(a) \propto \frac{\sqrt{B}}{a^3} \text{ for } a \ll \left(\frac{B}{A}\right)^{1/6}, \quad \rho_\varphi(a) \propto \sqrt{A} \text{ for } a \gg \left(\frac{B}{A}\right)^{1/6}. \quad (11.26)$$

At early times this gas behaves like pressureless dust, like CDM, at late times it behaves like the cosmological constant, causing accelerated expansion. The problem is, that when the dust effect disappears but CDM remains, this causes a strong ISW effect and loss of CMB power, thus the model is a poor fit to data (SNe, BAO, CMB). To remedy this one has complicated the model by adding a new parameter, β , to the barotropic equation of state:

$$p_\varphi = -A/\rho_\varphi^\beta. \quad (11.27)$$

Unfortunately, the fit then approaches the standard cosmological constant solution.

11.3 $f(R)$ Models

In Section 5.5 we have already met models of modified gravity which are candidates both for inflation and for the present accelerated expansion.

In Equation (5.85) we replaced the curvature scalar R by a general function $f(R, \varphi)$, and enlarged the dimensionality of space-time from 4 to n .

The Einstein–Hilbert action is not renormalizable, therefore standard general relativity cannot be properly quantized. However, renormalizability can be cured by adding higher order terms in curvature invariants. This is similar to the situation in quantum field theory where *renormalization* is some procedure to remove infinities in calculations. If the Lagrangian contains combinations of field operators of high

enough dimension in energy units, the counterterms required to cancel all divergences proliferate to infinite number, and, at first glance, the theory would seem to gain an infinite number of free parameters and therefore lose all predictive power, becoming scientifically worthless. Such theories are called nonrenormalizable.

The Standard Model of particle physics contains only renormalizable operators, but the interactions of general relativity become nonrenormalizable operators if one attempts to construct a field theory of quantum gravity in the most straightforward manner, suggesting that perturbation theory is useless in applications to quantum gravity.

The endeavor to remedy the renormalizability of the Einstein–Hilbert action implies adding functions of the Ricci scalar R which become important at late times and for small values of the curvature. This is required in order to avoid conflicts with constraints from the solar system or the Galaxy.

Consider an action of the form

$$S \propto \left[\int d^4x \sqrt{-g} [R + f(R)] + \int d^4x \sqrt{-g} \mathcal{L}_m(g_{\mu\nu}, \Psi) \right] \tag{11.28}$$

where $f(R)$ is an unspecified function of R . The matter Lagrangian \mathcal{L}_m is minimally coupled and, therefore the matter fields Ψ fall along the metric $g_{\mu\nu}$. By varying this action with respect to $g_{\mu\nu}$ one can obtain the field equations of the form

$$\tilde{G}_{\mu\nu} \equiv R_{\mu\nu} - \frac{1}{2}Rg_{\mu\nu} - Q_{\mu\nu} \tag{11.29}$$

where $\tilde{G}_{\mu\nu}$ contains the modifications to the geometry in terms of the functions $f(R)$, $f_R \equiv \partial f / \partial R$, $g_{\mu\nu}$. The $f(R)$ term leads to extra terms in the Einstein equation which is of second order in the metric $g_{\mu\nu}$. The constraint equation of GR becomes a third order equation, and the evolution equations now become fourth order differential equations which are difficult to analyze. A more tractable form of extended gravity can be obtained by using the *Palatini variation*, briefly discussed at the end of Chapter 5.

When taking the metric to be of the Robertson–Walker form [Equation (2.32)] the Friedmann equation becomes

$$H^2 + \frac{f}{6} - \frac{\ddot{a}}{a} f_R + H \dot{f}_R = \frac{\kappa^2 \rho}{3} \tag{11.30}$$

where $f \equiv f(R)$. The Raychauduri equation becomes

$$\frac{\ddot{a}}{a} - f_R H^2 + \frac{f}{6} + \frac{\dot{f}_R}{2} = -\frac{\kappa^2}{6}(\rho + 3p) \tag{11.31}$$

and the stress-energy tensor is replaced by

$$\tilde{T}_{\mu\nu} = \frac{T_{\mu\nu}}{f_R}. \tag{11.32}$$

Any $f(R)$ theory designed to achieve cosmic acceleration must satisfy $|f \ll R|$ and $|f_R| \ll 1$ at high curvature to be consistent with our knowledge of the high redshift universe. In order for f_R not to be tachyonic it must have a positive squared mass, that is, there must exist a stable high-curvature regime, such as a matter dominated

universe $f_{RR} > 0$ for $R \gg f_{RR}$. Further requirements are $1 + f_R > 0$ for all finite R (this prevents the graviton from becoming ghost-like), $f(R)/R \rightarrow 0$ and $f_R \rightarrow 0$ as $R \rightarrow \infty$ (GR must be recovered early before the BBN and the CMB), and f_R must be small at recent epochs [13]. The condition for an accelerated expansion is $a(t) \propto t^d$ with $d > 1$, suitable also for the exponential quintessence potential.

An example of a good potential is

$$f(R) = R \left[1 + \alpha \left(-\frac{R}{H_0^2} \right)^{\beta-1} \right] \quad (11.33)$$

with α and β constants, and which describes well the radiation-dominated era, the matter-dominated era, and the present accelerated era, but not the inflation. In the limit $R = 0$ one has flat geometry and standard GR.

Higher Order Invariants. Nothing forbids one to complement the Ricci scalar R with invariants of higher order in curvature. The advantage is that they introduce extra degrees of freedom, but they may come at a cost. They can lead to instabilities and conflicts with local tests of gravity. The invariants of lowest mass dimension are

$$P \equiv R_{\mu\nu}R^{\mu\nu} \quad Q \equiv R_{\alpha\beta\gamma\delta}R^{\alpha\beta\gamma\delta}. \quad (11.34)$$

which have actions similar to Equation (11.27),

$$S \propto \left[\int d^4x \sqrt{-g} [R + f(R, P, Q)] + \int d^4x \sqrt{-g} \mathcal{L}_m(g_{\mu\nu}, \Psi) \right]. \quad (11.35)$$

The equations of motion can again be found by varying this action with respect to $g_{\mu\nu}$.

The combination of invariants that seems most promising and is most studied is $R^2 - 4P + Q$, called the *Gauss-Bonnet term*.

11.4 Extra Dimensions

Already the generalized form of the Einstein–Hilbert action [Equation (5.85)] allowed for the possibility of a spacetime of $n > 4$ dimensions. Let us assume that the observable physics occurs on a four-dimensional brane embedded in a five-dimensional Minkowski bulk universe. Gravitation occurs everywhere in this bulk, but at cosmologically late times (now) it gets weaker on the brane, it leaks out into the bulk outside our brane, causing the present accelerated expansion.

The DGP model. A simple and well-studied model of modified gravity in five space-time dimensions is the *Dvali–Gabadadze–Porrati (DGP) braneworld model*. The model is characterized by a cross-over scale r_c such that gravity is a four-dimensional theory at scales $a \ll r_c H_0$ where matter behaves as pressureless dust. In the *self-accelerating* DGP branch, gravity ‘leaks out’ into the bulk when $a \approx r_c H_0$, and at scales $a \gg r_c H_0$ the model approaches the behavior of a cosmological constant. To explain the accelerated

expansion which is of recent date ($z \approx 0.7$ or $a \approx 0.6$), $r_c H_0$ must be of the order of unity. In the *self-decelerating* DGP branch, gravity ‘leaks in’ from the bulk at scales $a \approx r_c H_0$, in conflict with the observed dark energy acceleration. Note that the self-accelerating branch has a ghost, whereas the self-decelerating branch is ghost-free.

On the four-dimensional brane the action of gravity is proportional to M_{Pl}^2 whereas in the bulk it is proportional to the corresponding quantity in five dimensions, M_5^3 . The cross-over length is defined as

$$r_c = M_{\text{Pl}}^2 / 2M_5^3. \quad (11.36)$$

It is customary to associate a density parameter to this,

$$\Omega_{r_c} = (2r_c H_0)^{-2}, \quad (11.37)$$

such that $r_c H_0$ is a length scale (similar to a).

The Friedmann–Lemaître equation may be written

$$H^2 + \frac{k}{a^2} - \epsilon \frac{1}{r_c} \sqrt{H^2 + \frac{k}{a^2}} = \kappa \rho, \quad (11.38)$$

where $a = (1+z)^{-1}$, $\kappa = 8\pi G/3$, and ρ is the total cosmic fluid energy density with components ρ_m for baryonic and dark matter, and ρ_ϕ for whatever additional dark energy may be present.

Clearly the standard FLRW cosmology is recovered in the limit $r_c \rightarrow \infty$ or $H \ll r_c$. When $H \geq r_c$ the root term becomes important. In flat-space geometry $k = 0$, and at late times when $\rho \propto 1/a^3 \rightarrow 0$ and $H \rightarrow H_\infty$ this becomes a de Sitter acceleration, $a(t) \propto \exp t/r_c$.

The *self-accelerating branch* corresponds to $\epsilon = +1$, the *self-decelerating branch* to $\epsilon = -1$. In DGP geometry the continuity equations for ideal fluids have the same form as in FLRW geometry [Equation (11.28)].

Pressureless dust with $p = 0$ then evolves as $\rho_m(a) \propto a^{-3}$. The free parameters in the DGP model are Ω_{r_c} and $\Omega_m = \kappa \rho_m / H_0^2$. Note that there is no curvature term Ω_k since we have assumed flatness by setting $k = 0$ in Equation (11.33).

In the space of the parameters Ω_k , Ω_m and $\Omega_{r_c} \equiv 1/4r_c^2 H_0^2$ one can generalize the DGP model to

$$H^2 - \frac{k}{a^2} - r_c^{-2} \left(r_c \sqrt{H^2 - \frac{k}{a^2}} \right)^{2-n} = \kappa \rho, \quad (11.39)$$

where n defines a parametric family. Comparisons with data show that $n = 1$ and $n > 3$ give poor fits, $n = 2$ corresponds to the λ CDM concordance model, and $n = 3$ gives good fits.

The Chaplygin–DGP Model. Both the self-accelerating DGP model and the standard Chaplygin gas model have problems fitting present observational data, both cause too much acceleration. They have at least one parameter more than λ CDM, yet they fit data best in the limit where they reduce to λ CDM.

Because of the similarities in the asymptotic properties of the two-parametric DGP model and the two-parametric Chaplygin gas model. I have proposed [13] to combine the *self-decelerating* branch of the DGP model with the accelerating Chaplygin gas model. There are then four free parameters, Ω_{r_c} , Ω_m , A , and B , one of which shall be eliminated shortly.

We now choose the length scales in the two models, $r_c H_0$ and $(B/A)^{1/6}$, to be proportional by a factor x , so that

$$\left(\frac{B}{A}\right)^{1/6} = x r_c H_0 = \frac{x}{2\sqrt{\Omega_{r_c}}}. \quad (11.40)$$

The proportionality constant subsequently disappears because of a normalizing condition at $z = 0$. Then the model has only one parameter more than the standard Λ CDM model.

It is convenient to replace the parameters A and B in Equation (11.40) by two new parameters, $\Omega_A = H_0^{-2} \kappa \sqrt{A}$ and $x = 2\sqrt{\Omega_{r_c}} (B/A)^{1/6}$. The dark energy density can then be written

$$\rho_\phi(a) = H_0^2 \kappa^{-1} \Omega_A \sqrt{1 + x^6 (4\Omega_{r_c} a^2)^{-3}}. \quad (11.41)$$

Let us now return to Equation (11.38) and solve it for the expansion history $H(a)$. Substituting Ω_{r_c} from Equation (11.37), $\rho_\phi(a)$ from Equation (11.41), and using $\Omega_m = \Omega_m^0 a^{-3}$, it becomes

$$\frac{H(a)}{H_0} = -\sqrt{\Omega_{r_c}} + \left[\Omega_{r_c} + \Omega_m^0 a^{-3} + \Omega_{r_c}^0 a^{-4} + \Omega_A \sqrt{1 + x^6 (4\Omega_{r_c} a^2)^{-3}} \right]^{1/2}. \quad (11.42)$$

Note that Ω_{r_c} and Ω_A do not evolve with a , just like Ω_λ in the the λ CDM model. In the limit of small a this equation reduces to two terms which evolve as $a^{-3/2}$, somewhat similarly to dust with density parameter $\sqrt{\Omega_m^0 + \Omega_A x^3 (4\Omega_{r_c})^{-3/2}}$.

In the limit of large a , Equation (11.42) describes a de Sitter acceleration with a cosmological constant $\Omega_\lambda = -\sqrt{\Omega_{r_c}} + \sqrt{\Omega_{r_c} + \Omega_A}$.

A closer inspection of Equation (11.42) reveals that it is not properly normalized at $a = 1$ to $H(1)/H_0 = 1$, because the right-hand-side takes different values at different points in the space of the parameters Ω_m^0 , Ω_{r_c} , Ω_A , and x . This gives us a condition: at $a = 1$ we require that $H(1) = x H_0$ so that Equation (11.42) takes the form of a sixth order algebraic equation in the variable x

$$x = -\sqrt{\Omega_{r_c}} + \left[\Omega_{r_c} + \Omega_m^0 + \Omega_A \sqrt{1 + x^6 (4\Omega_{r_c})^{-3}} \right]^{1/2}. \quad (11.43)$$

This condition shows that x is a function $x = f(\Omega_m^0, \Omega_{r_c}, \Omega_A)$. Finding real, positive roots x and substituting them into Equation (11.42) would normalize the equation properly. The only problem is that the function cannot be expressed in closed form, so one has to resort to numerical iterations. The average value of x is found to be $x \approx 0.956$; it varies over the interesting part of the parameter space, but only by ≈ 0.002 .

This model [13] fits SNeIa data with the same goodness of fit as the the cosmological constant model, it also fits the CMB shift parameter R well, and notably, it offers a

genuine alternative to the cosmological constant model because it does not reduce to it in any limit of the parameter space.

The Gauss–Bonnet Model. Consider an action in a six-dimensional spacetime of the form

$$S \propto \left[\int d^6x \sqrt{-g} (R + \epsilon \mathcal{L}_{GB}) \right] \quad (11.44)$$

where ϵ is the *Gauss–Bonnet (GB)* parameter. Only $\epsilon > 0$ yields accelerating solutions. The metric consists of two pieces, the FRW metric and the GB metric,

$$ds^2 = c^2 dt^2 - a(t)^2 dl_{\text{FRW}}^2 - b(t)^2 [(dx^4)^2 + (dx^5)^2]. \quad (11.45)$$

The FRW Hubble flow and acceleration are the standard ones, $H = \dot{a}/a$ and \dot{H} , respectively. The GB Hubble flow is $h = \dot{b}/b$ and the GB acceleration is \dot{h} . The model has (at least) two interesting solutions:

- (i) In GB space there is decreasing deceleration, $\dot{h} < 0$, towards a stop when $\dot{h} = 0$. With time, the FRW behavior starts to dominate. In FRW space the expansion H and acceleration $\dot{H} > 0$ continue until a future singularity as in de Sitter expansion.
- (ii) H and h oscillate around future values of expansion and contraction. In FRW space there is almost no acceleration, $\dot{H} \approx 0$. In GB space \dot{h} oscillates between acceleration and deceleration, finally arriving at a future fixed point.

Lemaître–Tolman–Bondi Models. A way to explain the dimming of distant supernovae without dark energy may be an inhomogeneous universe described by the Lemaître–Tolman–Bondi solution (LTB) to Einstein’s equation. We could be located near the center of a low-density void which would distort our measurements of the age of the Universe, the CMB acoustic scale, and the Baryon Acoustic Oscillations (BAO).

The LTB model describes general radially symmetric spacetimes in four dimensions. The metric can be described by

$$ds^2 = -\alpha^2 dt^2 + X^2(r, t) dr^2 + A^2(r, t) d\Omega^2, \quad (11.46)$$

where $d\Omega^2 = d\theta^2 + \sin^2\theta d\varphi^2$, $A(r, t)$ and $X(r, t)$ are scale functions, and $\alpha(t, r) > 0$ is the *lapse function*.

Assuming a spherically symmetric matter source with baryonic+dark matter density ρ_M , dark energy density ρ_{DE} , negligible matter pressure $p_M = 0$ and dark energy pressure density p_{DE} , the stress-energy tensor is

$$T_{\nu}^{\mu} = p_{DE} g_{\nu}^{\mu} + (\rho_M + \rho_{DE} + p_{DE}) u^{\mu} u_{\nu}. \quad (11.47)$$

The $(0, r)$ components of Einstein’s equations, $G_r^0 = 0$, imply

$$\frac{\dot{k}(t, r)}{2[1 - k(t, r)]} + \frac{\alpha' \dot{A}}{\alpha A'} = 0 \quad (11.48)$$

with an arbitrary function $k(t, r)$ playing the rôle of the spatial curvature parameter. Here an overdot is designating ∂_t , and an apostrophe ∂_r .

The T_0^0 component of Einstein's equations gives the Friedmann–Lemaître equation in the LTB metric [Equation (11.46)]

$$\frac{H_T^2}{\alpha^2} + \frac{2H_T H_L}{\alpha^2} + \frac{k}{A^2} + \frac{k'(t, r)}{AA'} = 8\pi G (\rho_M + \rho_{DE}). \quad (11.49)$$

Here $H_T = \dot{A}/A$ is the transversal Hubble expansion, $H_L = \dot{A}'/A'$ the longitudinal Hubble expansion.

The T_r^r components of Einstein's equations give

$$2\dot{H}_T + 3H_T^2 + \frac{k}{A^2} - 8\pi G p_{DE} = 0. \quad (11.50)$$

Multiplying each term with $A^3 H_T$ this can be integrated over time to give

$$H_T^2 = \frac{F(r)}{A^3} - \frac{k(r)}{A^2} - \frac{8\pi G}{A^3} \int dt A^3 H_T p_{DE}. \quad (11.51)$$

Here $F(r)$ is an arbitrary time-independent function which arose as an integration constant.

We now assume that dark energy does not interact with matter. There are then separate continuity equations for ordinary matter and dark energy. The one for matter is

$$\dot{\rho}_M(r, t) + [2H_T(r, t) + H_L(r, t)] \rho_M(r, t) = 0, \quad (11.52)$$

and it can be integrated to give

$$\rho_M(r, t) = \frac{F(r)}{A^2 A'} \quad (11.53)$$

The continuity equation for dark energy is

$$\dot{\rho}_{DE}(r, t) + [2H_T(r, t) + H_L(r, t)] [\rho_{DE}(r, t) + p_{DE}(r, t)] = 0. \quad (11.54)$$

To proceed we should now specify several arbitrary functions, so we leave the subject here unfinished.

Problems

1. Derive Equation (5.76).
2. What should t_{eq} be for K-essence ρ_k to drop precisely to the magnitude of the present-day $\rho_\lambda \approx 2.9 \times 10^{-47} \text{ GeV}^4$?
3. Suppose that dark energy is described by an equation of state $w = -0.9$ which is constant in time. At what redshift did this dark energy density start to dominate over matter density? What was the radiation density at that time?
4. Derive this expression for the K-essence energy density:

$$\rho_k = [2XK'(X) - K(X)] + V(\varphi)$$

5. Derive the field equations for an action of the form in Equation (11.28) with $f(R) = R^2$. The ordinary matter part \mathcal{L}_m can be ignored.

References

- [1] Riess, A. G. *et al.* 1998 *Astronom. J.* **116**, 1009.
- [2] Perlmutter, S. *et al.* 1998 *Nature* **391**, 51; 1999 *Astrophys. J.* **517**, 565.
- [3] de Bernardis, F. *et al.* 2000 *Nature* **404**, 955.
- [4] Hanany, S. *et al.* 2000 *Astrophys. J. Lett.* **545**, L5.
- [5] Farook, O. and Ratra, B. 2013 *Astrophys. J. Lett.* **766**, L7.
- [6] Busca, N. G. *et al.* 2012 *Astron. Astrophys* **552**, A96.
- [7] Lima, J. A. S. and Trodden, M. 1996 *Phys. Rev. D* **53**, 4280.
- [8] Cunha, J. V., Lima, J. A. S. and Pires, N. 2002 *Astron. Astrophys.* **390**, 809.
- [9] Dimopoulos, K. and Valle, J. W. F. 2002 *Astroparticle Phys.* **18**, 287.
- [10] Steinhardt, P. J., Wang, L. and Zlatev, I. 1999 *Phys. Rev. D*, **59**, 123504.
- [11] Zlatev, I., Wang, L. and Steinhardt, P. J. 1999 *Phys. Rev. Lett.* **82**, 896.
- [12] Armendariz-Picon, C., Mukhanov, V. and Steinhardt, P. J. 2000 *Phys. Rev. Lett.* **85**, 4438; 2001 *Phys. Rev. D*, **63**, 103510.
- [13] Roos, M. 2008 *Phys. Letters B* **666**, 420.

12

Epilogue

We have now covered most of cosmology briefly, starting from the theoretical scenarios in Figure 5.1 and arriving at Figure 12.1, where the scenarios have been populated by data which permit a selection. Thus we know now that the cosmic expansion is accelerated since recently, since the time corresponding to redshift $z \approx 0.8$. But we still do not know what causes the acceleration—therefore the mysterious term ‘dark energy’, which represents 71.4% of the energy density of the Universe. In Chapter 11 we met a multitude, albeit not exhaustive set of candidates. The observational knowledge has improved greatly, but the solution remains unknown. The simplest candidate, the cosmological constant λ , begs an answer as to its origin and its value which has to be fine-tuned to within the 52nd decimal of zero (in units of $c = 1$). Determinations of the equation of state of dark energy, w , point to phantom models, but the cosmological constant value $w = -1$ is still close to the 1σ contour, see Figure 8.7 and the cover of this book.

The knowledge of the equally mysterious dark matter has gone from a disbelieved curiosity in 1933 to the important field of study today described in Chapter 8, much thanks to the enormous development in gravitational lensing. It represents 23.4% of the energy density of the Universe, yet its composition is unknown, except that it is cold. What has happened most recently is, that the long favored WIMP candidates have reached the useful search limit, and that the minimal supersymmetry model appears unlikely. In any case one should keep in mind, that the discovery of a new particle in a terrestrial accelerator does not prove that it is dark matter. That should be discovered at galactic scales.

Cosmic inflation is a field of vigorous study, but we are still relying on essentially the same model as 20 years ago, as described in Section 7.2. One should construct a renormalized field theory of quantum gravity to cope with the initial Big Bang singularity and inflation, but that is beyond the level of this monography. Models of inflation will be tested in the near future with the advent of more and better evidence for gravitational waves.

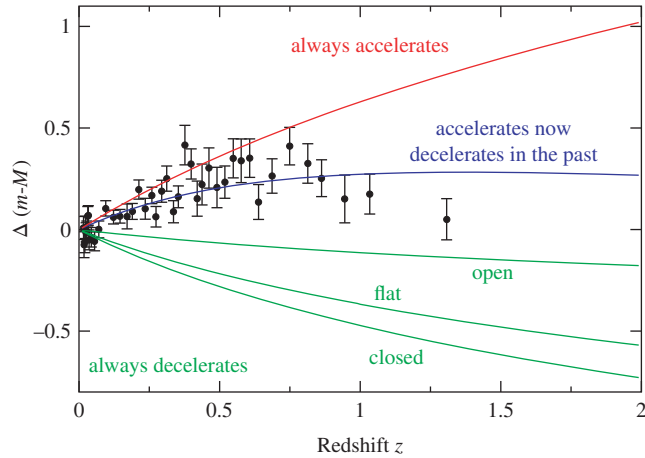


Figure 12.1 Evidence for transition from deceleration in the past to acceleration today. The difference $m-M$ is defined in Equation (2.60). The blue line shows a model that fits the data, where acceleration happens at late epochs in the history of the universe (i.e. starting a few billion years ago, and billions of years after the Big Bang). The plot uses binned data from the Union2 compilation [2]. With permission from the author. (See plate section for color version.)

Occasionally we have had reason to discuss a universe with extra dimensions: the cyclic universe in Section 7.4, some models of dark energy in Section 11.4. They remain candidates.

The future will bring great improvements in observational capacity, as we have witnessed ever since COBE, so cosmology will remain an active and fascinating field.

I thank the reader for her/his patience to come to this sentence!

References

- [1] Huterer, D. 2011 *Adventures in Cosmology*, ed. D. Goodstein, World Scientific Pub. Co., and preprint arXiv:1010.1162 [astro-ph.CO].
- [2] Amanullah, R. *et al.* 2010 *Astrophys. J.* **716**, 712.

Tables

Table A.1 Cosmic distances and dimensions

Distance to the Sun	8' 15" (light minutes)
Distance to the nearest star (α Centauri)	1.3 pc
Diameters of globular clusters	5–30 pc
Thickness of our Galaxy, the 'Milky Way'	0.3 kpc
Distance to our galactic center	8 kpc
Radius of our Galaxy, the 'Milky Way'	12.5 kpc
Distance to the nearest galaxy (Large Magellanic Cloud)	55 kpc
Distance to the Andromeda nebula (M31)	770 kpc
Size of galaxy groups	1–5 Mpc
Thickness of filament clusters	$5 h^{-1}$ Mpc
Distance to the Local Supercluster center (in Virgo)	17 Mpc
Distance to the 'Great Attractor'	$79 h^{-1}$ Mpc
Size of superclusters	$\gtrsim 50 h^{-1}$ Mpc
Size of large voids	$60 h^{-1}$ Mpc
Distance to the Coma cluster	$100 h^{-1}$ Mpc
Length of filament clusters	$100 h^{-1}$ Mpc
Size of the 'Sloan Great Wall'	$420 h^{-1}$ Mpc
Hubble radius	$3000 h^{-1}$ Mpc

Table A.2 Cosmological and astrophysical constants (from the Particle Data Group compilation [1])

Unit	Symbol	Value
Speed of light	c	$299\,792\,458\text{ m s}^{-1}$
Light year	ly	$0.3066\text{ pc} = 0.9461 \times 10^{16}\text{ m}$
Parsec	pc	$3.262\text{ ly} = 3.085\,678 \times 10^{16}\text{ m}$
Solar luminosity	L_{\odot}	$(3.846 \pm 0.008) \times 10^{26}\text{ J s}^{-1}$
Solar mass	M_{\odot}	$1.989 \times 10^{30}\text{ kg}$
Solar equatorial radius	R_{\odot}	$6.955 \times 10^8\text{ m}$
Hubble parameter	H_0	$100\text{ h km s}^{-1}\text{ Mpc}^{-1} = h/(9.778\,13\text{ Gyr})$
	h	0.696 ± 0.007
Newtonian constant	G	$6.673 \times 10^{-11}\text{ m}^3\text{ kg}^{-1}\text{ s}^{-2}$
Planck constant	\hbar	$6.582\,119 \times 10^{-22}\text{ MeV s}$
Planck mass	$M_{\text{P}} = \sqrt{\hbar c/G}$	$1.221 \times 10^{19}\text{ GeV } c^2$
Planck time	$t_{\text{P}} = \sqrt{\hbar G/c^5}$	$5.31 \times 10^{-44}\text{ s}$
Boltzmann constant	k	$8.617\,34 \times 10^{-5}\text{ eV K}^{-1}$
Stefan–Boltzmann constant	$a = \pi^2 k^4/15\hbar^3 c^3$	$4.7222 \times 10^{-3}\text{ MeV m}^{-3}\text{ K}^{-4}$
Critical density of the Universe	$\rho_{\text{C}} = 3H_0^2/8\pi G$	$2.775 \times 10^{11}\text{ h}^2 M_{\odot}\text{ Mpc}^{-3}$ $= 10.538\text{ h}^2\text{ GeV m}^{-3}$

Table A.3 Electromagnetic radiation

Type	Wavelength [m]	Energy [eV]	Energy density ¹ [eV m ⁻³]
Radio	> 1	$< 10^{-6}$	≈ 0.05
Microwave	$1\text{--}5 \times 10^{-3}$	$10^{-6}\text{--}2 \times 10^{-3}$	3×10^5
Infrared	$2 \times 10^{-3}\text{--}7 \times 10^{-7}$	$10^{-3}\text{--}1.8$?
Optical	$(7\text{--}4) \times 10^{-7}$	$1.8\text{--}3.1$	$\approx 2 \times 10^3$
Ultraviolet	$4 \times 10^{-7}\text{--}10^{-8}$	$3.1\text{--}100$?
X-rays	$10^{-8}\text{--}10^{-12}$	$100\text{--}10^6$	75
γ rays	$< 10^{-12}$	$> 10^6$	25

¹From M. S. Longair [2].**Table A.4** Particle masses¹

Particle	MeV units	K units
γ	0	0
$\langle v_i \rangle$	$< 0.23 \times 10^{-6}$	$< 2.7 \times 10^3$
e^{\pm}	0.511	5.93×10^9
μ^{\pm}	105.658	1.226×10^{12}
π^0	134.977	1.566×10^{12}
π^{\pm}	139.570	1.620×10^{12}
p	938.272	1.089×10^{13}
n	939.565	1.090×10^{13}
τ^{\pm}	1777	2.062×10^{13}
W^{\pm}	80423	9.333×10^{14}
Z	91188	1.058×10^{15}
H^0	125.3 ± 0.6	1.454×10^{15}

¹From the Particle Data Group compilation [1].

Table A.5 Particle degrees of freedom in the ultrarelativistic limit

Particle	Particle type	n_{spin}	n_{anti}	g
γ	Vector boson	2	1	2
ν_e, ν_μ, ν_τ	Fermion (lepton)	1 ¹	2 ²	$\frac{7}{4}$
e^-, μ^-, τ^-	Fermion (lepton)	2	2	$\frac{7}{2}$
π^\pm, π^0	Boson (meson)	1	1	1
p, n	Fermion (baryon)	2	2	$\frac{7}{2}$
W^\pm, Z	Vector boson	3	1	3

¹ $n_{\text{spin}} = 2$, but the right-handed neutrinos are inert below the electroweak symmetry breaking.

² $n_{\text{anti}} = 1$ if the neutrinos are their own antiparticles.

Table A.6 Present properties of the Universe

Unit	Symbol	Value
Age	t_0	(13.73 ± 0.05) Gyr
Mass	M_U	$\approx 10^{22} M_\odot$
CMB radiation temperature	T_0	2.725 ± 0.001 K
Cosmic neutrino temperature	T_ν	1.949 ± 0.001 K
Radiation energy density	$\epsilon_{r,0}$	2.606×10^5 eV m ⁻³
Radiation density parameter	Ω_r	$2.473 \times h^{-2} \times 10^{-5}$
Entropy density	s	2.890×10^9 m ⁻³
CMB photon number density	N_γ	4.11×10^8 photons m ⁻³
Cosmological constant	$ \lambda $	$1.3 \times 10^{-52} c^2$ m ⁻²
Schwarzschild radius	$r_{c,\text{Universe}}$	$\gtrsim 11$ Gpc
Baryon to photon ratio	η	$(6.06 \pm 0.08) \times 10^{-10}$
Vacuum density parameter	Ω_k	-0.003 ± 0.003
Baryon density parameter (for $\Omega_0 = 1$)	Ω_b	0.052
Matter density parameter (for $\Omega_0 = 1$)	Ω_m	0.286 ± 0.008
Deceleration parameter	q_0	-0.53 ± 0.02
Number of neutrino families	N_ν	3.30 ± 0.27

Table A.7 Net baryon number change ΔB and branching fraction BR for leptoquark X decays

i	Channel i	ΔB_i	BR _{i}
1	$X \rightarrow \bar{u} + \bar{u}$	$-\frac{2}{3}$	r
2	$X \rightarrow e^- + d$	$+\frac{1}{3}$	$1 - r$
3	$\bar{X} \rightarrow u + u$	$+\frac{2}{3}$	\bar{r}
4	$\bar{X} \rightarrow e^+ + d$	$-\frac{1}{3}$	$1 - \bar{r}$

References

[1] Beringer, J. *et al.* 2012 *Phys. Rev. D* **86**, Part I, 010001.
 [2] Longair, M. S. 1995 *The Deep Universe* (ed. A. R. Sandage, R. G. Kron, and M. S. Longair), pp. 317, Springer.

Index

Note: Figures are indicated by *italic page numbers*, Tables by **bold page numbers**

- Abell 2744 cluster, 212
- absolute luminosity, 9, 42, 44
- absolute magnitude, 42, 190
- absolute space, 5, 7, 30
- absorption lines, 28, 133, 140
- ACT-CLJ0102 – 4915 cluster, 212
- action principle, 58, 59
- active galactic nuclei (AGN), 17, 106
- adiabatic expansion, 86, 115, 120
- adiabatic fluctuations, 169, 182, 231
- affine connections, 56, 107
- age of the Universe, 14, 92, 195, 244, **259**
- Alpher, Ralph, 176
- Andromeda nebula, 5, 6, 12, 177, 201
 - distance to, 257
- angular size distance, 43
- anisotropy
 - quadrupole, 75, 182, 185, 188
 - sources of, 182
- annihilation, 101
 - baryon – anti-baryon, 142
 - electron – positron, 101, 121, 129
 - leptoquark boson, 145
 - monopole – anti-monopole, 156
 - pion, 127
 - WIMP, 216
- anthropic principle, 151, 169
- anthropocentric view, 2
- antibaryons, 17
- anti-bias, 209
- anti-de Sitter universe, 93
- antigravity, 146, 236
- antineutrinos, 17
- antiparticles, 121
- anti-protons, 121
- apparent luminosity, 9
- apparent magnitude, 42
- astrophysical constants [listed], **258**
- autocorrelation function
 - mass, 207, 225
 - temperature, 180
- axino, 216
- axion, 216
- B-balls, 214
- Baby Bullet, 212
- BAO (baryonic acoustic oscillations), 14
- baryon, 135, 138, 143, 168, 194, 216, 231
 - mirror, 214
- baryon – anti-baryon asymmetry, 142, 168
- baryon number, 122
 - per unit entropy, 146
- baryon number density, 135, 142
- baryon-to-photon ratio, 138, 142, 194, **259**
- baryonic matter, 17
 - dark, 213
- baryosynthesis, 142, 146
- Bekenstein, J., 101
- Bekenstein – Hawking formula, 101
- beta decay, 137

- bias, 209
- BICEP2 telescope, 185, 189, 193
- Big Bang, 88, 89, 99, 112, 154, 177, 241
 - nucleosynthesis, 16, 17, 111, 134
- Big Crunch, 89
- Big Rip, 245
- binary star systems, 18, 67, 104
- binding energy, 132
- BIT (backwards in time) interpretation, 148
- black holes, 4, 54, 76, 96, 124, 214, 226
 - creation of, 103
 - firewall paradox, 100
 - Kerr – Newmann black hole, 99
 - merging of, 101, 102
 - observations, 104
 - Reissner – Nordström black hole, 99
 - Schwarzschild black hole, 94, 96
- blackbody spectrum, 116, 177
- blue compact dwarf (BCD) galaxies, 139
- blueshift, 30
- Bohr orbits, 27
- Boltzmann, Ludwig, 116
 - see also* Stefan – Boltzmann law
- Boltzmann constant, **258**
- Bose, Satyendranath, 123
- Bose distribution, 120
- Bose – Einstein condensate (BEC), 215
- bosons, 120, 123
 - Higgs, 122, 123, 144
 - scalar, 239
 - vector, 121, 122, 123, 144, 157
- bottom – top scenario, 218
- bouncing universe, 151, 169, 190
- bound state, 132, 135
- Bradley, James, 2
- brane, 171, 214, 245, 248
- Brans – Dicke theory, 107
- Brightest Cluster Galaxies (BCGs), 16
- brightness
 - apparent, 42
 - surface, 9, 44, 73
- brown dwarfs, 213
- bubble universe, 151, 162, 167
- Bullet cluster, 211

- CDM paradigm, 218, 220
 - see also* cold dark matter
- Cepheids, 5, 44, 73
- Chandrasekhar mass, 18
- chaotic inflation, 151, 159, 165
- Chaplygin gas model, 246

- charge conjugation, 145
- charged current, 127
- Chéseaux, Jean-Philippe Loys de, 8
- classical mechanics, 19, 224
- closed gravitational system, 8, 22
- cluster flattening parameter, 210, 221
- clusters, dark matter in, 208
- CMB, 14, 175
 - polarization of, 185
 - reionization of, 191, 231
 - temperature anisotropies, 180, 223
 - see also* cosmic microwave background radiation
- COBE (Cosmic Background Explorer), 175, 178, 180, 184
- cold dark matter, 215
- collapsed stars, 213
- collisional dissipation, 231
- Coma cluster, 14, 44, 208
 - distance to, 257
- comoving coordinates, 34, 35, 36
- comoving distance, 37
- comoving frame, 36
- Compton scattering, 121, 129, 185
- Compton wavelength, 101, 112
- concordance model, 81, 151, 193, 235, 249
- conformal time, 37
- conformal transformation, 107
- consensus inflation, 151, 158, 171
- conservation laws, 122
- consistency test, 14
- contraction operation, 57
- contravariant vector, 54
- cooling plasma, 125
- Copernican principle, 3, 22, 58
- Copernicus, Nicolaus, 2
- cored profile, 201
- cosmic
 - ensorship, 99
 - inflation, 112, 151
 - rays, 17, 143
 - scale factor, 13
 - shear field, 73, 74
 - strings, 156, 182
 - structures, 223
 - time, 36
- cosmic distances and dimensions [listed], **257**
- cosmic microwave background (CMB) radiation, 14, 175
 - discovery of, 176

- temperature, 176, **259**
see also CMB
- cosmochronometers, 15, 190
- cosmological coincidence problem, 236
- cosmological constant, 84, 85, 93, 159, 235, **259**
decaying, 238
- cosmological constant problem, 236
- cosmological constants [listed], **258**
- cosmological principle, 3, 6
- cosmological redshift, 13
- Coulomb force, 115
- coupling constants, 115
- covariance principle, 54
- covariant derivative, 56
- covariant vector, 54
- CP violation, 144, 145
- CPT symmetry, 147
- critical curves, 74
- critical density [of the Universe], 20, **258**
- critical temperature, 113, 161
- curvature parameter, 35
- curvature perturbations, 169, 182, 183
- curved manifold, 35
- curved space – time, 30
- cusped profile, 201
- cyclic universe, 151, 169, 190

- dark energy, 19, 168, 172, 218, 230, 235
- dark matter, 19, 169, 199
 - baryonic, 213
 - candidates, 213
 - cold, 215
 - distribution, 216
 - halo density profiles, 200
 - hot, 217
 - warm, 218
- DASI (Degree Angular Scale Interferometer), 175, 189
- de Sitter, Willem, 6
- de Sitter cosmology, 91, 93, 162, 243
- de Sitter metric, 93, 94
- deceleration, 82
- deceleration parameter, 39, 76, 195, **259**
- decoupling, 134, 195
 - electron, 135
 - neutrino, 130
- deflation, 241
- degeneracy pressure, 18, 123
- degrees of freedom, 123, **259**
effective, 124

- density fluctuations, 223, 227
- density parameters, 20, 83, 85, 87, 132, 191, 193
values listed, **259**
- density perturbations, 164
- deuterium, 135, 140
bottleneck, 136
- deuteron(s), 135
photodisintegration of, 135
- DGP (Dvali – Gabadadze – Porrati) model, 248
self-accelerating branch, 248
self-decelerating branch, 248, 249
- Dicke, Robert, 177
see also Brans – Dicke theory
- dipole anisotropy, 180
- Dirac, Paul A. M., 156
see also Fermi – Dirac statistics
- distance ladder, 43
- DLSCJ0916.2+2951 cluster, 213
- DMR (Differential Microwave Radiometer), 180, 184
- domain walls, 156
- Doppler blueshift, 30
- Doppler peak, 184
- Doppler redshift, 30
- dust, interstellar, 5, 9, 17, 87, 90, 213
- dwarf spheroidal galaxies, 206

- Eddington, Arthur S., 70
- Einstein, Albert, 4, 6
- Einstein – de Sitter universe, 83, 89, 157, 190
- Einstein
 - frame, 107
 - gravity, 59
 - ring, 71, 72
 - universe, 84, 93
- Einstein – Hilbert action, 58, 106, 160, 246, 247
in Jordan frame, 107
- Einstein’s mass – energy relation, 55
- Einstein’s theory of general relativity, 6, 49
- Einstein’s theory of special relativity, 25
- ‘El Gordo’ cluster, 212
- electromagnetic cross-section, 129
- electromagnetic interactions, 115, 122, 129, 133
- electromagnetic radiation, **258**
- electron, 17
- electron degeneracy pressure, 18, 123
- electron – positron reactions, 101, 121, 129

- electroweak interactions, 121
- elliptical galaxies, 204
- endothermic reactions, 134, 135
- energy conservation law, 86, 115, 120
- energy density, 117, 179
- energy effect, 42
- energy – momentum conservation, 85, 239
- energy – momentum tensor, 60, 169
- entropy, 101, 125, 160
- entropy conservation law, 86, 116, 130
- entropy density, 179, **259**
- equation of state, 86, 159, 195
- equilibrium theory, 133
- equivalence principle, 50
- eternal inflation, 162
- Euclidean space, 30
- Eulerian equations, 224
- event horizon, 39, 96, 99, 103
- Evershed, John, 67
- exothermic reactions, 129, 134, 135
- expansion rate, 128, 129
- expansion time, 190
- expansion velocities, 12
- extended gravity models, 106

- falling photons, 52
- false vacuum, 161
- Fermi, Enrico, 123
- Fermi coupling, 128
- Fermi – Dirac statistics, 123
- Fermi distribution, 120
- fermion, 120, 123
- fermion number, 123
- 'Fingers of God', 209, 221
- FIRAS (Far Infrared Absolute Spectrophotometer), 178
- first law of thermodynamics, 116
- first-order phase transition, 112
- flatness problem, 151, 157
- FLRW (Friedmann – Lemaître – Robertson – Walker) concordance model, 81, 151, 193, 235
- fluid dynamics, 224
- flux, 73
- Friedmann, Alexandr, 6
- Friedmann's equations, 81, 89, 91, 93
- Friedmann – Lemaître cosmologies, 81, 84, 235
- fundamental observer, 36
- fundamental plane, 45
- fusion reactions, 135

- Galaxy *see* Milky Way Galaxy
- galaxy clusters, 218, 226
 - merging of, 211
- galaxy formation, 207
- galaxy groups, 201
 - size, **257**
- Galilean equivalence principle, 51
- Galilei, Galileo, 2
- gamma ray bursts (GRBs), 17, 106
- gamma rays, 143
- Gamow, Georg, 176
- Gamow penetration factor, 138
- gauge, 164
- gauge interactions, 214
- gauge problem, 227
- gauge transformation, 227
- Gauss, Carl Friedrich, 33
- Gauss – Bonnet term, 248
- Gaussian curvature, 33, 35
- general covariance, 49, 55, 58
- general relativity, 6, 27, 49, 111
 - classical tests, 65
- geodesic, 30
 - in Minkowski space-time, 32
 - on sphere, 32
- geodesic equation, 57
- global positioning system, 67
- globular clusters, 5, 16, 44, 232
 - diameters, **257**
- gluons, 144
- gold, relativistic energies, 27
- GPS, 67
- 'graceful exit', 160, 228
- grand unified theory (GUT), 114, 143, 144
 - phase transition, 155
- gravitating mass, 19, 49
- gravitational
 - birefringence, 54
 - collapse, 18, 103, 106, 230, 231
 - fields, timekeeping in, 67
 - lensing, 66, 69
 - potential, 61
 - radiation, 68
 - repulsion, 85
- gravitational waves, 61, 68, 74, 165, 185
 - detection of, 75, 76
 - sources, 76
- gravitino, 218
- graviton, 75
- gravity
 - Chaplygin – DGP, 249

- Chaplygin gas, 246
 DGP, 248
 Einstein, 59
 extended, 106
 $f(R)$, 246
 Gauss – Bonnet, 251
 k-essence, 244
 Lemaitre – Tolman – Bondi, 251
 phantom fields, 245
 quintessence, 241
 tachyon fields, 245
 ‘Great Attractor’, 41, 221
 distance to, **257**
 great circles, 32
 Guth, Alan, 160
- Halley, Edmund, 2
 Hawking, Stephen, 99
 Hawking radiation, 100, 101, 165
 Hawking temperature, 102
 HDM, 217
 helium, 16, 137
 ions, 137
 Helmholtz, Hermann von, 120
 Herman, Robert, 176
 Herschel, William, 5
 Hertzprung – Russell relation, 42
 hierarchical scenarios [cold dark matter model], 218
 hierarchy problem, 115
 Higgs boson, 122, 123, 144
 Higgs field, 155
 Higgs-like potential, 159
 Higgsino, 215
 Hilbert, David, 58
 see also Einstein – Hilbert action
 homogeneity assumption, 2
 horizon problem, 151, 153
 HST (Hubble Space Telescope), 14, 44, 211, 212
 Hubble, Edwin P., 5, 6
 Hubble constant, 13
 Hubble flow, 12, 14, 38, 209, 224
 Hubble’s law, 2, 6, 11, 12, 14, 18, 41, 66
 Hubble parameter, 12, 125, **258**
 Hubble radius, 12, 39, **257**
 Hubble time, 12
 Hulse, R. A., 67
 Hydra – Centaurus supercluster, 30, 41, 180, 221
 hydrogen
- atom, 121, 132
 burning of, 15, 43, 140
 clouds, 139, 191, 199, 231
 ions, 137
 hypothesis testing, 14
- ideal fluid, 59
 inertial frames, 6
 inertial mass, 20, 49
 inflation
 chaotic, 151, 159, 165
 confirmation of, 193
 consensus, 151, 158, 171
 cosmic, 151
 eternal, 162
 slow-roll, 151, 158
 inflaton field, 158, 239, 243
 information loss paradox, 100
 infrared light, 43, 73
 interaction
 electromagnetic, 115, 122, 129, 133
 weak, 121
 see also gravitational repulsion
 intergalactic medium (IGM), 17, 140, 143
 interstellar medium, 140, 143
 intracluster medium (ICM), 201, 208
 isocurvature fluctuations, 169, 182, 183
 isothermal perturbations, 169
 isotropy assumption, 2
- Jeans, James, 230
 Jeans instability, 230
 Jeans mass, 228, 230
 Jeans wavelength, 230, 231
 Jordan frame, 107
 Jupiters, 213
- Kant, Immanuel, 3
 Kapteyn, Jacobus C., 204
 Kelvin, Lord [William Thomson], 9, 10
 Kepler, Johannes, 2
 Kerr – Newmann black hole, 99
 kination, 240, 241
 KK (Kaluza Klein) states, 214
 Klein – Gordon equation, 158, 239
- Lagrange point, 51
 Lagrangian density, 107
 Lambert, Johann Heinrich, 5
 Landau damping, 216
 Laplace, Pierre Simon de, 4
 lapse function, 251

- Large Magellanic Cloud (LMC), 18, 44, 73
 distance to, **257**
- Laser Interferometer Space Antenna *see*
 LISA
- last scattering, 134
- last scattering surface (LSS), 133, 134, 153,
 195
- Le Verrier, Urban, 66
- Legendre polynomials, 181
- Leibnitz, Gottfried Wilhelm von, 3
- Lemaître, Georges, 82
- Lemaître cosmology, *91*
see also Friedmann – Lemaître
 cosmologies
- Lemaître – Tolman – Bondi model, 251
- lens
 caustics, 74
 convergence, 74
 shear, 74
- lensing
 strong, 71
 weak, 69, 71, 74
- lepton number, 122
- leptoquark thermodynamics, 144
- light *see* speed of light
- light cone, 28, 29
- light element abundance, 139
- light-like separation, 28
- LIGO-type interferometer, 77
- Lindblad, Bertil, 6
- Linde's Bubble Universe, 162, 167, 190
- line element, 26
- linear transformation, 26
- LISA detector, 77, 78
- lithium, 139
- local galaxy group, 30, 41, 180, 201
- Local Supercluster *see* LSC
- lookback time, 83, 92, 244
- loops, 156
- Lorentz, Hendrik Antoon, 26
- Lorentz invariance, 45
- Lorentz transformations, 26
- LSC (Local Supercluster), 3, 41, 180, 202,
 209, 221
 distance to, **257**
- luminosity, 9, 42, 44
- luminosity distance, 42, 190
- Lyman limit, 140
- Lyth bound, 165
- M31 *see* Andromeda nebula
- Mach, Ernst, 5
- Mach's principle, 49
- MACHO (Massive Compact Halo Object),
 213
- magnetic monopoles, 156
- magnitude
 absolute, 42, 190
 apparent, 42
- magnitude – redshift relation, 190
- main-sequence stars, 42
- manifold, 26
 curved, 35
 higher-dimensional, 55
- mass autocorrelation function, 207, 225, 227
- mass density contrast, 224
- mass of the Universe, **259**
- matter domination, 87, 88, 89, 117, 170
- Maxwell, James Clerk, 120
- Maxwell – Boltzmann distribution, 120
- Maxwell – Lorentz equations, 115
- mean free path, 10
- Mercury, 45, 66
- metals, meaning of term [to astronomers],
 139
- metric
 flat, 36
 Gauss – Bonnet, 251
 Lemaître – Tolman – Bondi, 251
 Minkowski, 28, 31, 32
 Pythagorean, 30
 Robertson – Walker, 36, 56, 81, 107, 247
 Schwarzschild, 95, 100
- metric tensor, 30
- Michell, John, 4
- microlensing, 73
- Milky Way Galaxy, 2 – 6, 15, 16, 18, 19, 43,
 44, 73, 140, 143, 180, 201, 204
 dimensions, **257**
- Minkowski, Hermann, 27
- Minkowski space-time, 27, 31
- mirror matter, 214
- monopole problem, 151, 156
- multipole analysis, 180, 188
- muons, 26, 128
- naked singularity, 99
- neutralino, 215
- neutrino(s), 17
 clouds, 218
 decoupling of, 130
 families, 138, 194, **259**

- number density, 131, 179
- oscillation of, 122, 146
- sterile, 215, 218
- temperature contribution, 124, 131
- neutron degeneracy pressure, 18, 123
- neutron star, 18, 19, 67, 68, 103, 123, 126
- neutron-to-proton ratio, 134, 138
- Newton, Isaac, 2, 5
- Newton's first law of motion, 6
- Newton's law of gravitation, 2, 19, 22, 49
- Newton's second law of motion, 49, 55, 56
- Newtonian constant, 19, 45, **258**
- Newtonian cosmology, 2, 5
- Newtonian mechanics, 19, 22
- nonrelativistic particles, 118
- nuclear fusion, 135
- null separation, 28
- number density, 116, 127, 131, 179
- object horizon, 39
 - see also* particle horizon
- Olbers, Wilhelm, 8
- Olbers' Paradox, 2, 8
- Oort, Jan Hendrik, 6
- open gravitational system, 8, 21
- optical depth, 191
- oscillating universe, 151, 169, 190
- Palatini variation, 107, 247
- parallax distance, 42
- parallel axiom, 33
- parameter estimation, 14, 189
- parity operator, 145
- parity transformation, 145
- Parker bound, 156
- parsec [unit], 6, **258**
- particle horizon, 39, 94, 103, 152
- particle masses [listed], **258**
- Pauli, Wolfgang, 123
- Pauli exclusion force, 123
- peculiar velocity, 14, 209
- Penrose, Roger, 99
- Penzias, Arno, 176
- perihelion, 66
- perturbations, 164, 169
- phase transitions, 112
- photino, 215
- photon
 - blackbody spectrum, 116, 177
 - diffusion of, 230
 - entropy density, 130, 146
 - number density, 116, 135, 179
 - reheating of, 129
 - spin states, 187
 - superluminal, 54
 - virtual, 121
- Pioneer anomaly, 45
- pions, reactions and decay, 127, 128
- Planck, Max, 53, 116
- Planck constant, 53, **258**
- Planck mass, 112, **258**
- Planck space mission, 175
- Planck time, 112, **258**
- Poisson's equation, 61, 224
- polarization
 - circularly polarized light, 186
 - CMB, 185
 - fluctuations, 185
 - linearly polarized light, 186
- positronium, 121
- positrons, 17, 121
- power spectrum, 181, 191, 225
- powers, 181
- preferred direction of time, 161, 170
- present observable universe, 28
- pressure
 - matter, 87
 - radiation, 87, 224, 231
 - vacuum energy, 87, 239
- primeval asymmetry generation, 143
- primordial density fluctuations, 207
- primordial plasma, 112, 124
- Proctor, Richard Anthoy, 5
- proper distance, 37, 40
- proper time, 26
- protons, 121, 122
- pulsars, binary, 67
- Q-balls, 214
- QED (quantum electrodynamics), 121
- quadrupole anisotropy, 75, 182, 185, 188
- quantum fluctuations, 166
- quantum mechanics, 53
- quark matter, 19, 126
- quark star, 103
- quasar, 17
- quintessence, 239, 241
 - tracking, 242
- R parity, 115
- radiation domination, 87, 88, 170, 177
- radiation energy density, 179, **259**

- radiation intensity, 178, 188
- radiation pressure, 87, 224, 231
- radio signal delay, 67
- radius of the Universe, 90, 165
- rank, 55
- Raychauduri equation, 82, 159, 247
- reaction rate, 128
- recession velocities, 12, 66
- recombination, 134
- recombination era, 129, 132, 153
- recombination radiation, 129
- recombination redshift, 134, 176
- recombination time, 134, 176
- red giant, 18
- redshift, 28
 - cosmological, 13, 29, 40
 - distance, 42
 - Doppler, 30
 - and luminosity distance, 42
 - and proper distance, 40
- reheating, 161, 168
 - of photon, 129
- reionization, 134, 191
- Reissner – Nordström black hole, 99
- relativistic particles, 118
- relativity
 - general, 6, 27, 49, 111
 - special, 25
- relic ^4He abundance, 138
- renormalization, 246
- Ricci scalar, 58, 59, 82, 90, 107
- Ricci tensor, 57, 82
- Riemann, Bernhard, 4
- Riemann tensor, 57, 82
- Riemann's zeta-function, 116
- Robertson, Howard, 36
- Robertson – Walker metric, 36, 56, 81, 107, 247
- rotating galaxies, 3
- RR Lyrae stars, 43

- Sachs – Wolfe effect, 183
- Saha equation, 133
- Sakharov oscillations, 184
- scalar spectral index, 169
- scalar-tensor theories, 107
- scale factor, 29, 30
- Schwarzschild, Karl, 95
- Schwarzschild black hole, 94, 96
- Schwarzschild metric, 95, 100
- Schwarzschild radius, 95, 97, 98, 99, 104, **259**
- second cosmic velocity, 8
- second law of thermodynamics, 116, 120, 130
- second-order phase transition, 112
- Shapiro, I. L., 67
- Shapley, Harlow, 5
- Shen, Yang, 4
- Silk, J., 231
- Silk damping, 231
- silver, relativistic energies, 27
- slicing, 164
- Sloan Great Wall, 218, 219, **257**
- slow-rolling conditions, 159, 243
- snowballs, 213
- solar constant(s), 148, **258**
- Solar System, 1 – 5, 15, 19, 30, 65, 143, 180, 204
- solitons, 214
- sound, velocity of, 229
- space-time distance, 26
- spacelike separation, 28
- sparticles, 115
- special relativity, 25
 - tests, 45
- speed of light, 12, 13, 25, 54, **258**
 - variable, 45
- Spherical Collapse Model, 208
- spin states, 187
 - longitudinal, 123
 - transversal, 123
- spiral galaxies, 203, 204
- standard candle, 18, 43
- standard model, 114
- star formation, 15, 139, 207, 221
- starlight deflection, 66, 71
- statistics, 14
- Stefan, Josef, 116
- Stefan – Boltzmann constant, **258**
- Stefan – Boltzmann law, 116
- sterile neutrinos, 215, 218
- Stokes parameters, 186, 187
- stress – energy tensor, 60, 159, 238, 251
- strong equivalence principle (SEP), 51, 69
- strong lensing, 71
- structure formation, 228
 - time, 232
- structure size, 232

- structure stimulation, 220
- subconstituent models, 114
- Sunyaev – Zel’dovich Effect (SZE), 191, 212, 231, 232
- superclusters, 14, 41, 180, 209, 218
 - size, **257**
- superluminal photons, 54
- supernovae, 5, 14, 15, 18, 76, 139
- supersymmetry (SUSY), 115, 144, 215
- surface-brightness fluctuations (SBFs), 44
- synchrotron, 27

- tachyons, 54, 245
- Taylor, J. H., 68
- technicolor forces, 115
- temperature
 - anisotropies, CMB, 180, 223
 - autocorrelation function, 180
 - critical, 113, 161
 - fluctuations, 180, 185
 - multipole components, 182
 - scale, relation to timescale, 125
- tensor, 30, 54
- tensor field, 75
- tensor spectral index, 169
- thermal conductivity, 230
- thermal death, 120
- thermal equilibrium, 116
- thermal history of the Universe, 111
- thermodynamics
 - first law, 116
 - second law, 116, 120, 130
- Thomson scattering, 129, 175, 185, 187, 231
- threading, 164
- tidal effect, 51
- time dilation, 26
- time reversal, 146
- timelike separation, 28
- timescale, 194
- timescale test, 190
- Tolman – Oppenheimer – Volkoff limit, 103
- Tolman test, 44
- top – bottom scenario, 218
- topological defects, 156
- tracking quintessence, 242
- trigonometrical parallax, 42
- tritium, 136
- triton, 136
- Tully – Fisher relation, 44

- turnover time, 89
- twin paradox, 67
- two-point correlation function, 226

- ultra-compact dwarf galaxies (UCDs), 207
- ultra-faint dwarf galaxies (UFDs), 207
- Universal Rotation Curve, 204
- universe
 - anti-de Sitter, 93
 - bouncing/cyclic/oscillating, 151, 169, 190
 - closed, 20, 37, 89
 - contracting, 8, 13, 22, 89
 - de Sitter, 93, 94, 162, 243
 - Einstein, 84, 93
 - Einstein – de Sitter, 83, 89, 157, 190
 - expanding, 8, 12, 13, 20, 21, 89
 - Friedmann, 94
 - Friedmann – Lemaître, 81, 84
 - inflationary, 160
 - Newtonian, 19
 - open, 20, 38
 - steady state theory, 89
 - thermal history, 111

- vacuum energy, 85, 87, 152, 168, 193
- vacuum energy density, 93
- vacuum energy pressure, 87, 239
- vector bosons, 121, 122, 123, 144, 157
- Virgo supercluster, 14, 41, 44, 76, 202
- virial equilibrium, 232
- virial theorem, 200
 - applications, 208
- virially bound systems, 200
- virtual particles, 101
- virtual photons, 121
- viscous fluid approximation, 224
- visibility limit, 218

- Walker, Arthur, 36
 - see also* Robertson – Walker metric
- warm dark matter, 218
- wavenumber, 181
- WDM, 218
- weak equivalence principle (WEP), 51
- weak field limit, 61
- weak interaction, 121
- weak lensing, 69, 71, 74
- Weyl, Hermann, 37
- Wheeler, John A., 96

white dwarfs, 18, 45, 123

white hole, 99

Wilson, Robert, 176

WIMPs (weakly interacting massive particles), 215, 216

WMAP (Wilkinson Microwave Anisotropy Probe), 175, 189

world line, 28, 32

wormhole, 98

Wright, Thomas, 3

Zel'dovich, Yakov B., 191

see also Sunyaev – Zel'dovich Effect

Zino, 215

Zwicky, Fritz, 73, 208

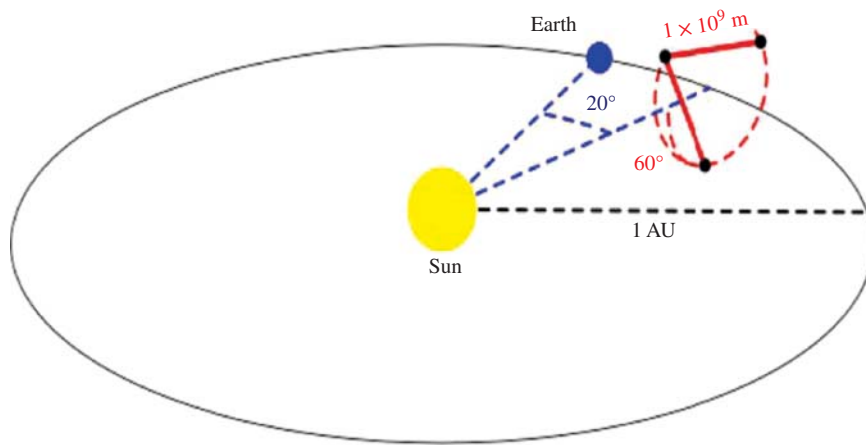


Figure 4.4 eLISA orbits in the Solar System [4]. Reproduced with permission of Pau Amaro-Seoane and the LISA consortium.

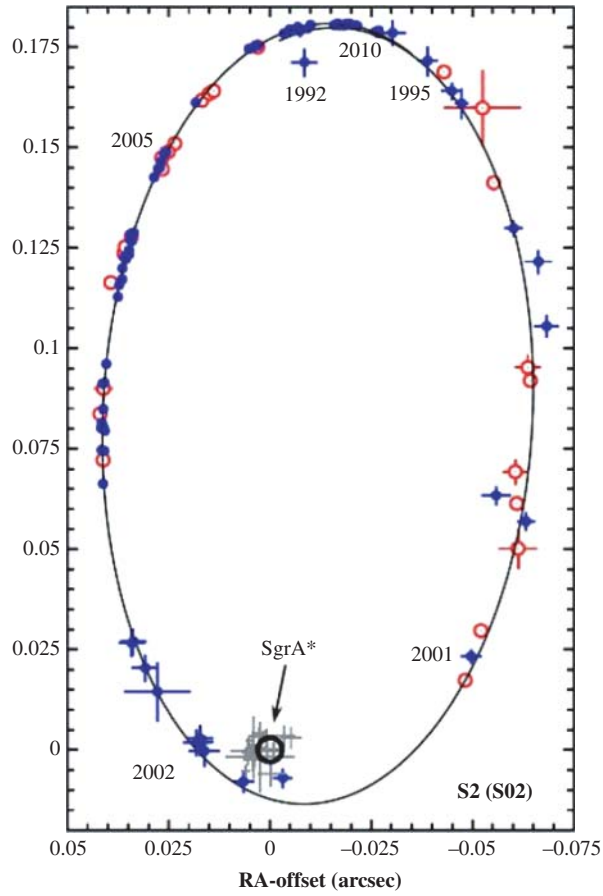


Figure 5.6 Orbit of the star S2 moving around Sgr A* [10]. Copyright 2010 by the American Physical Society.

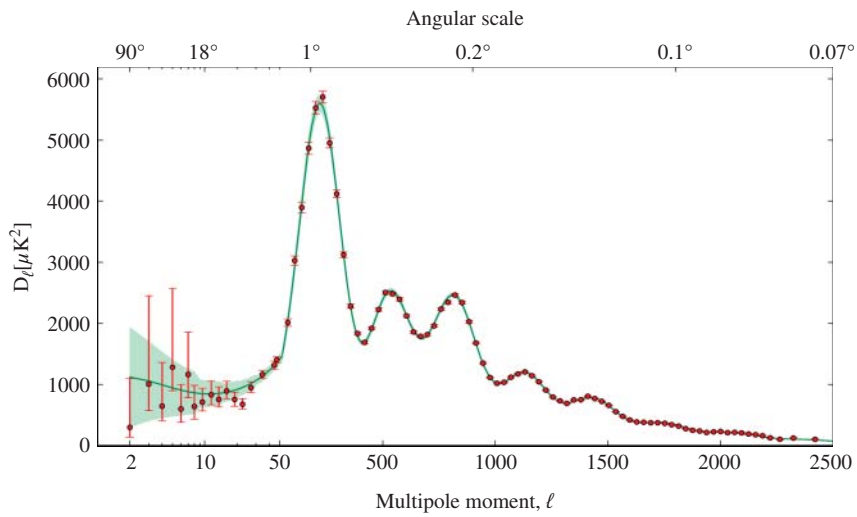


Figure 8.3 The best-fit power spectra of CMB temperature (T) fluctuations as a function of angular scale (top x axis) and multipole moment (bottom x axis) [6]. Reproduced from the freely accessible Planck Legacy Archive with permission of Jan Tauber, European Space Agency.

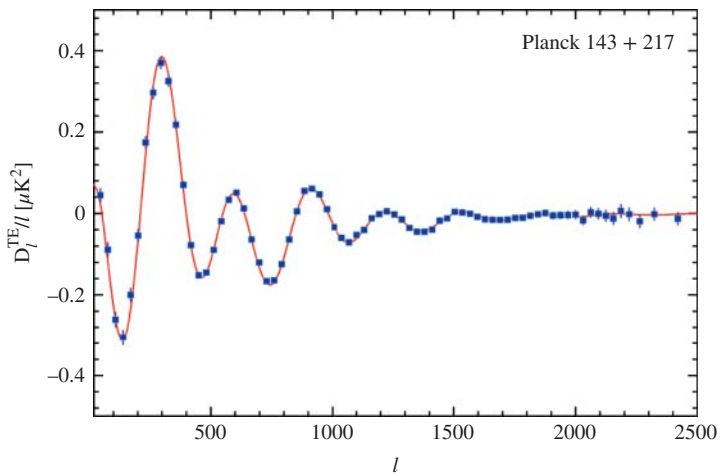


Figure 8.4 The temperature–E-polarization cross-power spectrum as a function of angular scale (top x axis) and multipole moment (bottom x axis) [6]. Reproduced from the freely accessible Planck Legacy Archive with permission of George Efstathiou, Kavli Institute for Cosmology, University of Cambridge.

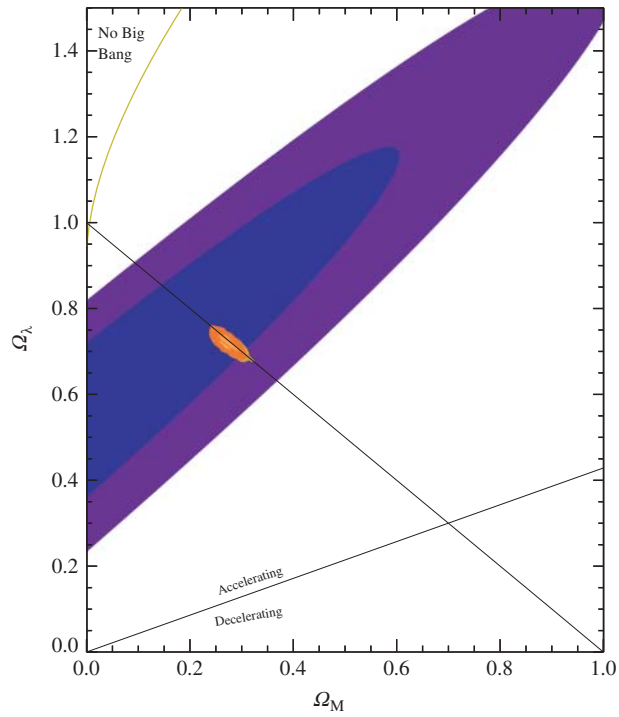


Figure 8.5 The 1σ and 2σ cosmological constraints in Ω_m , Ω_λ space using PanSTARRS1 supernova data [13], Planck CMB data [6], BAO [14] and H_0 data [6] with statistical and systematic errors propagated. Reproduced with permission of Armin Rest for the PanSTARRS1 Collaboration.

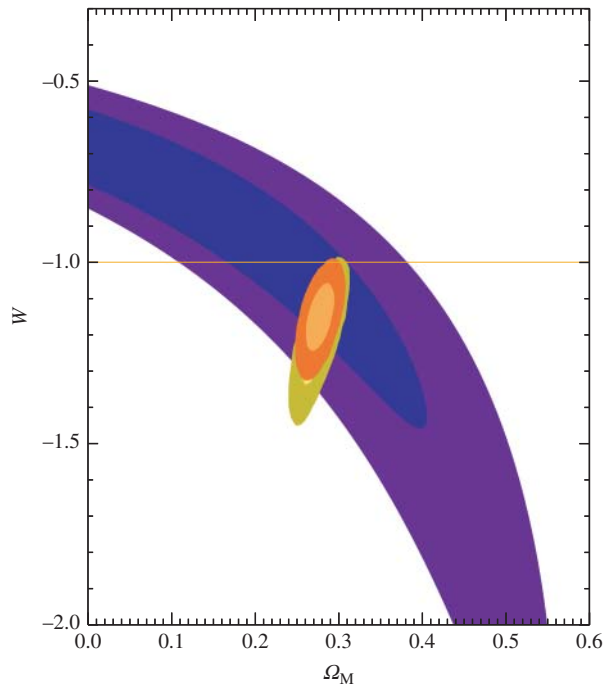


Figure 8.6 The 1σ and 2σ cosmological constraints in Ω_m , W space using PanSTARRS1 supernova data [13], Planck CMB data [6], BAO [14] and H_0 data [6] with statistical and systematic errors propagated. Reproduced with permission of Armin Rest for the PanSTARRS1 Collaboration.

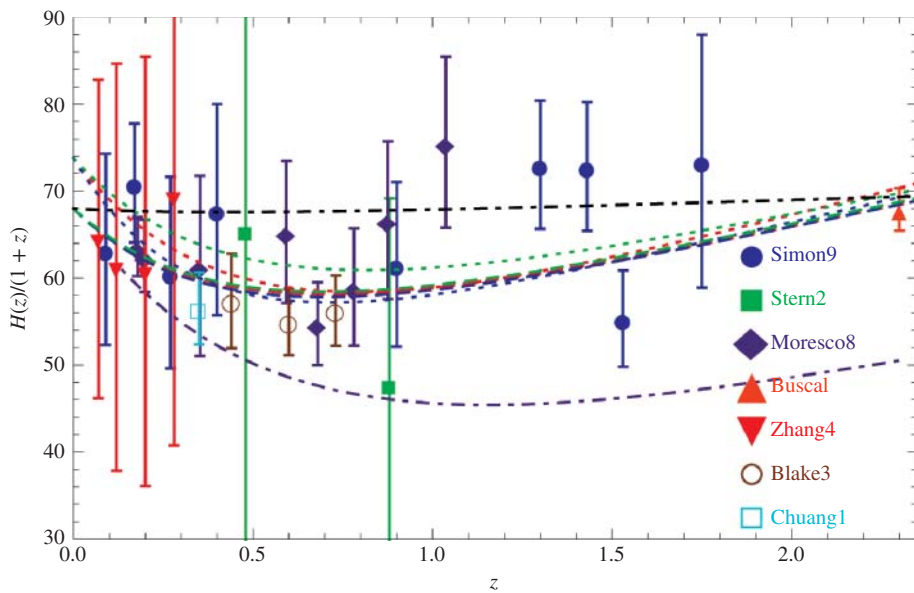


Figure 11.1 Evidence for transition from deceleration in the past to acceleration today. From reference [6]. From Farook, O. and Ratra, B., Hubble parameter measurement constraints on the cosmological deceleration-acceleration transition redshift, *Astrophys. J. Lett.*, **766**, L7, published 4 March 2013. © AAS. Reproduced with permission.

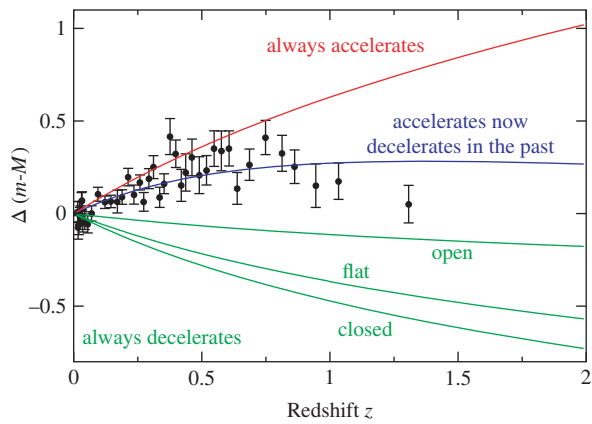


Figure 12.1 Evidence for transition from deceleration in the past to acceleration today. The difference $m-M$ is defined in Equation (2.60). The blue line shows a model that fits the data, where acceleration happens at late epochs in the history of the universe (i.e. starting a few billion years ago, and billions of years after the Big Bang). The plot uses binned data from the Union2 compilation [2]. With permission from the author.

WILEY END USER LICENSE AGREEMENT

Go to www.wiley.com/go/eula to access Wiley's ebook EULA.